

# Answering Conjunctive Queries and FO+MOD Queries under Updates

Dissertation  
zur Erlangung des akademischen Grades

doctor rerum naturalium (Dr. rer. nat.)

im Fach Informatik

eingereicht an der

Mathematisch-Naturwissenschaftlichen Fakultät der Humboldt-Universität zu Berlin

von

M.Sc. Jens Keppeler

Präsidentin der Humboldt-Universität zu Berlin

Prof. Dr.-Ing. Dr. Sabine Kunst

Dekan der Mathematisch-Naturwissenschaftlichen Fakultät:

Prof. Dr. Elmar Kulke

Gutachter/innen :      1. Prof. Dr. Nicole Schweikardt, Humboldt-Universität zu Berlin  
                                 2. Prof. Dr. Stefan Kratsch, Humboldt-Universität zu Berlin  
                                 3. Prof. Dr. Wim Martens, Universität Bayreuth

Tag der mündlichen Prüfung: 25. Mai 2020



# Abstract

This thesis investigates the query evaluation problem for fixed queries over fully dynamic databases, where tuples can be inserted or deleted. The task is to design a dynamic algorithm that immediately reports the new result of a fixed query after every database update.

In particular, the goal is to construct a data structure that allows to support the following scenario: after every database update, the data structure can be updated in constant time such that afterwards we are able to test within constant time for a given tuple whether or not it belongs to the query result (the testing routine), to output the number of tuples in the query result (the counting routine), to enumerate all tuples in the new query result (the enumeration routine), and to enumerate the difference between the old and the new query result with constant delay (the difference routine). The preprocessing time needed to build the data structure is linear in the size of the database.

In the first part of this thesis, conjunctive queries on arbitrary relational databases are considered. The notion of *q-hierarchical* conjunctive queries is introduced and it is shown that the result of each such query on a dynamic database can be maintained efficiently in the sense described above. Moreover, this notion is extended to aggregate queries for which we can also maintain the query result under updates. Furthermore, it is shown that the preparation of learning a polynomial regression function can be done in constant time if the training data are taken (and maintained under updates) from the query result of a q-hierarchical query. For the testing problem, the notion of t-hierarchical conjunctive queries, a more expressive query language, is considered and it is shown that such queries can be maintained efficiently and allow to test whether or not a given tuple belongs to the result set. It turns out that if one allows logarithmic update time in the size of the database (and  $n \log(n)$  preprocessing time as well), then the data structure for q-hierarchical queries additionally supports the following routine: upon input of a natural number  $j$ , output the  $j$ th tuple that will be enumerated by the enumeration routine. Furthermore, the notion of q-hierarchical and t-hierarchical conjunctive queries is extended to unions of conjunctive queries (UCQs) and it is shown that there is a data structure that can be maintained under updates, supports the enumeration routine and the testing routine, and a subset of q-hierarchical UCQs is considered for which one can output the  $j$ th solution of an enumeration.

In the second part of this thesis, queries in first-order logic (FO) and its extension with modulo-counting quantifiers (FO+MOD) are considered, and it is shown that they can be efficiently evaluated under updates, provided that the dynamic database does not exceed a certain degree bound, and the counting, testing, enumeration and difference routines are supported.



# Zusammenfassung

In dieser Arbeit wird das dynamische Auswertungsproblem für feste Anfragen über dynamischen Datenbanken betrachtet, bei denen Tupel hinzugefügt oder gelöscht werden können. Die Aufgabe besteht darin, einen dynamischen Algorithmus zu konstruieren, welcher unmittelbar nachdem die Datenbank aktualisiert wurde, eine Datenstruktur aktualisiert, die das Resultat ändert.

Darüber hinaus soll die Datenstruktur in konstanter Zeit aktualisiert werden und folgende Routinen unterstützen. Nach jeder Datenbankaktualisierung kann die Datenstruktur in konstanter Zeit angepasst werden, so dass anschließend in konstanter Zeit getestet werden kann, ob ein Tupel zur Ausgabemenge gehört (die Test-Routine), die Anzahl der Tupel in der Ausgabemenge in konstanter Zeit ausgegeben werden kann (die Anzahl-Routine), die Tupel aus der Ausgabemenge mit konstanter Taktung aufgezählt werden kann (die Aufzähl-Routine), und der Unterschied zwischen der neuen und der alten Ausgabemenge mit konstanter Taktung aufgezählt werden kann (die Unterschied-Routine). Die Vorverarbeitungszeit um die Datenstruktur aufzubauen benötigt lineare Zeit in der Größe der Datenbank.

Im ersten Teil dieser Arbeit werden konjunktive Anfragen auf beliebigen relationalen Datenbanken betrachtet. Die Idee der *q-hierarchischen* Anfragen wird eingeführt und es wird gezeigt, dass das Resultat für jede q-hierarchische Anfrage auf dynamischen Datenbanken effizient ausgewertet werden können in dem oben beschriebenen Szenario. Darüber hinaus wird dieses Szenario auf q-hierarchische Anfragen mit Aggregaten erweitert und es wird gezeigt, dass das gleiche Szenario ebenfalls effizient gelöst werden kann. Außerdem wird gezeigt, dass das Lernen von polynomiellen Regressionsfunktionen in konstanter Zeit vorbereitet werden kann, falls die Trainingsdaten aus dem Anfrageergebnis einer q-hierarchischen Anfragen kommen und insbesondere in konstanter Zeit aktualisiert werden können, wenn sich die Datenbank ändert. Für die Test-Routine betrachten wir die t-hierarchischen Anfragen, die ausdrucksstärker als q-hierarchische Anfragen sind, für die wir die Test-Routine auf dynamische Datenbanken auswerten können. Es stellt sich heraus, dass wenn man logarithmische Zeit in den Aktualisierungen für q-hierarchische Anfragen zulässt (und insbesondere  $n \log(n)$  Vorverarbeitungszeit erhält), man auch folgende Routine realisieren kann: Bei Eingabe einer natürlichen Zahl  $j$ , gib das  $j$ te Tupel von dem Aufzählalgorithmus aus. Außerdem, werden die Ideen von q-hierarchischen und t-hierarchischen konjunktive Anfragen auf Vereinigungen konjunktive Anfragen (UCQs) erweitert und es wird gezeigt, dass es eine Datenstruktur für diese Anfragen auf dynamische Datenbanken gibt, die mit konstanter Aktualisierungszeit, die Aufzähl-Routine und die Test-Routine unterstützt. Darüber hinaus wird das Problem für die Ausgabe des  $j$ ten Tupel auf eine Teilmenge für q-hierarchische Anfragen betrachtet.

Im zweiten Teil der Arbeit werden Anfragen, die Formeln der Logik erster Stufe (FO)

und deren Erweiterung mit Modulo-Zähl Quantoren (FO+MOD) sind, betrachtet, und es wird gezeigt, dass diese effizient unter Aktualisierungen ausgewertet können, wobei die dynamische Datenbank die Gradschranke nicht überschreitet, und bei der Auswertung die Zähl-, Test-, Aufzähl- und die Unterschied-Routine unterstützt werden.

# Contents

<b>1. Introduction</b>	<b>1</b>
1.1. The author's contribution of the published articles . . . . .	4
1.2. Related work . . . . .	4
1.3. Recommended reading order . . . . .	7
<b>2. Preliminaries</b>	<b>9</b>
<b>1. Conjunctive Queries under Updates</b>	<b>17</b>
<b>3. Main Theorems and Organisation of Part I</b>	<b>19</b>
3.1. An introduction to conjunctive queries . . . . .	19
3.2. Main theorems and organisation . . . . .	20
<b>4. Answering q-hierarchical conjunctive queries under updates</b>	<b>29</b>
4.1. The fundamental data structure . . . . .	29
4.2. Preprocessing and updating the data structure . . . . .	34
4.3. Answering Boolean q-hierarchical queries . . . . .	40
4.4. Testing q-hierarchical queries . . . . .	41
4.5. The decomposition lemma . . . . .	42
4.6. Enumerating the query result with constant delay . . . . .	44
4.7. Enumerating the difference . . . . .	48
4.8. Outputting the next smallest tuple . . . . .	63
<b>5. Testing t-hierarchical conjunctive queries under updates</b>	<b>69</b>
<b>6. Answering q-hierarchical conjunctive queries with aggregates under updates</b>	<b>73</b>
6.1. Examples for queries with aggregates . . . . .	73
6.2. A model to evaluate aggregations . . . . .	76
6.3. Syntax and semantics of aggregation expressions . . . . .	80
6.4. How to compute the q-tree of a q-hierarchical conjunctive query with aggregates . . . . .	83
6.5. Counting the number of tuples in $Q(D)$ . . . . .	86
6.6. The data structure for queries with aggregates . . . . .	88
6.7. A reduction to q-hierarchical queries without aggregates . . . . .	94

<b>7. An application concerning learning polynomial regression models over q-hierarchical queries</b>	<b>101</b>
7.1. Linear regression . . . . .	103
7.2. Polynomial regression . . . . .	108
<b>8. Outputting the jth solution for q-hierarchical conjunctive queries under updates</b>	<b>115</b>
8.1. Output the jth solution . . . . .	116
8.2. Proof of Lemma 8.2 . . . . .	123
8.3. jth reverse . . . . .	128
8.4. Proof for Theorem 3.9 . . . . .	129
8.5. Lower bounds . . . . .	131
<b>9. Answering unions of conjunctive queries under updates</b>	<b>133</b>
9.1. Enumerating and testing UCQs . . . . .	133
9.2. Reporting the jth solution . . . . .	136
 <b>II. Answering FO+MOD Queries under Updates on Bounded Degree Databases</b>	 <b>147</b>
<b>10. Further preliminaries and main results in Part II</b>	<b>149</b>
<b>11. Answering Boolean FO+MOD Queries Under Updates</b>	<b>155</b>
<b>12. Answering non-Boolean FO+MOD Queries Under Updates</b>	<b>159</b>
12.1. Technical Lemmas on Types and Spheres Useful for Handling Non-Boolean Queries . . . . .	159
12.2. Testing Non-Boolean FO+MOD Queries Under Updates . . . . .	162
12.3. Representing Databases by Coloured Graphs . . . . .	165
12.4. Counting Results of FO+MOD Queries Under Updates . . . . .	170
12.5. Enumerating Results of FO+MOD Queries Under Updates . . . . .	172
12.6. Refining the enumeration routine . . . . .	178
12.7. Enumerating the Difference . . . . .	183
 <b>13. Conclusion</b>	 <b>187</b>
<b>Bibliography</b>	<b>189</b>



# 1. Introduction

In this thesis we consider the algorithmic task of evaluating a query on fully dynamic databases, where tuples can be inserted or deleted. More generally, the aim is to design a data structure representing the query and the database, which can be updated every time a tuple will be inserted into or removed from the database and it solves at least one of following problems efficiently:

- answer Boolean queries (*answering problem*),
- enumerate the tuples in the result set without repetition with constant delay (*enumeration problem*),
- on input of a tuple, test if the tuple belongs to the result set (*testing problem*),
- enumerate the tuples that were added to / removed from the result set after the last update step (*diff-and-report problem*),
- enumerate the tuples that are added in / removed from the result set after some update steps (*difference problem*),
- output the number of tuples in the result set (*counting problem*),
- prepare the learning of the parameters for a polynomial regression function where the training set is taken from the query result (*learning problem*) and
- given a natural number  $j$ , output the  $j$ th tuple in the enumeration, if it exists ( *$j$ th tuple problem*).

We consider finite relational databases over a possibly infinite domain.

In the beginning, we obtain an initial database  $D_0$  and a query  $Q$ . In the *preprocessing phase* we construct the data structure that represents the database  $D_0$  and  $Q$ . Whenever we change the database, i.e., a tuple is being inserted or deleted, we modify the data structure in an *update step* such that it represents the updated database. The *update time* is the time needed to process an update step. In order to be efficient, our aim is to show that the update steps take constant or logarithmic time in the size of the database and in particular, the update time is way smaller than the time needed to recompute the entire query result.

This thesis is divided into two parts.

The first part's aim is to give algorithms for a subclass of conjunctive queries (CQ for short), unions of conjunctive queries (UCQ for short), and conjunctive queries

## 1. Introduction

with aggregates, all of which can be efficiently maintained under updates. It turns out that for the answering problem, the enumeration problem, the difference problem and the counting problem, the class of *q-hierarchical* queries are suitable for solving these problems under updates. The notion of *q-hierarchical* queries are strongly related to the *hierarchical* property that was introduced by Dalvi and Suciu in [34] and already played a central role for efficient query evaluation in various contexts (see Chapter 3 for a definition). The following is shown: after a linear time preprocessing phase (here we consider data complexity, i.e., linear time means linear in the size of the database) the *q-hierarchical* query can be maintained under updates and can solve the answering problem, the enumerating problem, the counting problem and the difference problem (see Theorem 3.3). This means that after every update the data structure can still solve these problems. It turns out that the testing problem can be done efficiently under updates for a superclass of *q-hierarchical* queries after a linear time preprocessing. We call these conjunctive queries *t-hierarchical* (see Definition 3.4 for a definition and Theorem 3.6 for the result).

Furthermore, we also consider extensions of conjunctive queries, queries with aggregates and unions of conjunctive queries. We define a class of conjunctive queries with aggregates, such that we can maintain queries (by solving the answering, enumerating, testing and the counting problem) under updates after a linear time preprocessing phase (see Theorem 3.7). We show that every *q-hierarchical* query with aggregates can be reduced to *q-hierarchical* conjunctive queries without aggregates. In particular, it follows that if a problem is tractable under updates for *q-hierarchical* conjunctive queries, then it is also tractable for *q-hierarchical* CQs with aggregates if we can maintain the corresponding aggregate functions in the query under updates in sufficient time. That means, for example, we only obtain constant update time, if we can update the values of the aggregation function for the result set in constant time.

As mentioned before, we consider the problem of learning a polynomial regression function, a problem which comes from the topic of *supervised machine learning*. The main idea of this topic is to find relations or patterns between data of a given training set. After using an algorithm that learns a model (“learn” means in this context to compute the model) we can use the model to obtain approximate solutions of other inputs. In this thesis, we consider the model of a polynomial regression, which is a polynomial function that minimizes the error between the solutions in the training set and the value of the function. Such a learning algorithm has two passes, where in the first pass precomputations are done for the second pass. In the second pass, numerical methods were used to approximate the parameters (for more details see Chapter 7). In this thesis we show how to efficiently maintain the computations that are done in the first pass under updates where the training set is the query result (see Theorem 3.8).

For a nice survey to machine learning see [54, 43].

For the task of solving the *j*th problem, we consider logarithmic update time rather than constant time. We show that after an  $n \log(n)$  time preprocessing phase (here  $n$  is the size of the database), we can output for a given natural number  $j$  the *j*th tuple in the result of a *q-hierarchical* CQ with or without aggregates on a database  $D$  under logarithmic updates (if the aggregates can be computed in logarithmic time or faster)

(see Theorem 3.9).

For unions of conjunctive queries we obtain results for the testing problem for a union of t-hierarchical conjunctive queries (this is a UCQ where every CQ in the query is t-hierarchical) and the enumeration problem for q-hierarchical UCQ (a UCQ where every CQ in the query is q-hierarchical) under updates. Furthermore, we consider the problem for reporting the  $j$ th tuple for UCQ. We identify a class of UCQs for which we can maintain this problem efficiently under updates. We call these queries strongly exhaustively q-hierarchical UCQs (see Theorem 3.12).

In the second part of this thesis, we consider queries in first-order logic (FO) and its extension with modulo-counting quantifiers (FO+MOD) and show that they can be efficiently evaluated under updates, provided that the dynamic database does not exceed a certain degree bound. We obtain in the second part the following result: let  $Q$  be a  $k$ -ary FO+MOD-query and  $d$  a degree bound on the database.  $Q$  and  $d$  are assumed to be fixed, i.e., in contrast to the database they do not change. With  $\|Q\|$  ( $\|D\|$ ) we denote the size of a query  $Q$  (database  $D$ , resp.). On input of an initial database  $D_0$ , we can construct in a linear time  $f(\|Q\|, d)\|D\|$  preprocessing phase a data structure that can be updated in time  $f(\|Q\|, d)$  and allows to

- immediately answer  $Q$  on  $D$  if  $Q$  is a Boolean query,
- test for a given tuple whether it belongs to the result set in time  $O(k^2)$ ,
- immediately output the number of tuples in the result set  $Q(D)$ ,
- enumerate the tuples in the result set  $Q(D)$  with delay  $O(k^2)$  and
- enumerate the tuples in  $Q(D^-) \setminus Q(D^+)$  and  $Q(D^+) \setminus Q(D^-)$  with delay  $O(k^2)$ , where  $D^-$  and  $D^+$  denote the database before and after performing the update operation, respectively.

The function  $f(Q, d)$  stands for a function of the form

$$f(Q, d) = 2^{d^{2^{O(\|Q\|)}}}.$$

In contrast to the difference enumeration in the thesis' first part, we enumerate here the difference immediately after exactly one update step.

The dynamic query evaluation algorithm crucially relies on the locality of FO+MOD and, in particular, an effective Hanf normal form for FO+MOD on databases of bounded degree recently obtained by Heimberg, Kuske and Schweikardt [55]. The organization of the second part is given in Chapter 10.

## 1. Introduction

### 1.1. The author's contribution of the published articles

The articles [15],[17],[19], [16],[18],[21] and [20] were published with co-authors, namely Christoph Berkholz and Nicole Schweikardt. The author contacted his co-authors to reaffirm that there are no disputes over the ownership of the material presented in this thesis.

[15] is a preprint of [16]. The author's main contribution in [16] lies on the upper bounds that are specified in Chapter 1,2,3,4,6 and 7 in [16]. Moreover, there are extensions of these results that are presented in this thesis.

[19] is a preprint of [21]. The author's essential contribution are the upper bound for testing t-hierarchical conjunctive queries and t-hierarchical unions of conjunctive queries and the upper bounds for enumerating q-hierarchical unions of conjunctive queries. These results can be found in Chapter 1,2,3 and 4 in [21]. Moreover, the queries in [16] did not consider constants whereas in [21] the result was lifted up to q-hierarchical queries with constants. In this thesis, q-hierarchical queries were introduced with constants and the results are shown with considering constants.

In Chapter 3 the main theorems of the first part are given together with a description which parts of the main theorem are results from [16] and [21] and which parts are new results.

[17] is a preprint of [18] and published in a journal [20]. This work was made in a close cooperation with the co-authors and the contribution of the authors is to be recognized equally. Part II is closely based on the journal version [20].

### 1.2. Related work

There is a lot of work known that studies the complexity of query evaluation in the static setting. Since there is a huge amount of such work, surprisingly little is known about evaluating queries under updates. The task of answering queries against a dynamic database has been studied under the name of *incremental view maintenance* (see e.g. [53, 30, 66, 68, 81]).

There is a framework, introduced by Patnaik and Immerman [83], called the *dynamic descriptive complexity* framework, which focuses on the expressive power of first-order logic on dynamic databases (see [88] for a survey). Updating this approach may take polynomial time (even in the case of conjunctive queries [96]). This setting is too expensive in the area for dynamic algorithms.

Some work in computation complexity of query evaluation under updates has been done by Björklund, Gelade and Martens [22] for XPath evaluation and for MSO queries on trees by Balmin, Papakonstantinou and Vianu [11], by Losemann and Martens [74], by Niewerth and Segoufin [80] and by Amarilli et al. [4].

In [59], the enumeration and testing problem under updates has been studied for q-hierarchical and (more general) acyclic CQs in a setting that is very similar to the thesis' first part setting. The *Dynamic Constant-delay Linear Representations* (DCLR) of [59] are data structures that use at most linear update time and solve the enumeration problem and the testing problem with constant delay and constant

testing time.

In case of the static setting, a lot research has been done. Below, there is an overview of known results.

**Complexity of Boolean Queries.** The complexity of answering Boolean conjunctive queries on a static database is fairly well understood. For every fixed database schema  $\sigma$ , extending a result of [51], Grohe [47] gave a tight characterisation of the *tractable* CQs under the complexity theoretic assumption  $\text{FPT} \neq \text{W}[1]$ : If we are given a Boolean CQ  $Q$  of size  $\|Q\|$  and a  $\sigma$ -database  $D$  of size  $\|D\|$ , then  $Q$  can be answered against  $D$  in time  $f(\|Q\|) \cdot \|D\|^{O(1)}$  for some computable function  $f$  if and only if the homomorphic core of  $Q$  has bounded treewidth. Marx [76] extended this classification to the case where the schema is part of the input.

**Counting Complexity.** For computing the number of output tuples of a given join query (i.e., a quantifier-free CQ) over a fixed schema  $\sigma$ , a characterisation was proven by Dalmau and Jonsson [33]: assuming  $\text{FPT} \neq \#\text{W}[1]$ , the output size  $|Q(D)|$  of a join query  $Q$  evaluated on a  $\sigma$ -database  $D$  of size  $\|D\|$  can be computed in time  $f(\|Q\|) \cdot \|D\|^{O(1)}$  if and only if  $Q$  has bounded treewidth. The result has recently been extended to all conjunctive queries over a fixed schema by Chen and Mengel [28]. Structural properties that make the counting problem for CQs tractable in the case where the schema is part of the input have been identified in [37, 46]. The counting problem, i.e., output the number in the result set on input of a database and a query, has been studied for conjunctive queries [28] and for existential positive formulas [29]. In [33] Dalmau and Jonsson showed that the problem of counting the number of homomorphisms in a class  $\mathcal{C}$  of structures to a given arbitrary structure is solvable if and only if all structures in  $\mathcal{C}$  have bounded treewidth.

**Join Evaluation.** When the entire result of a non-Boolean query has to be computed, the evaluation problem cannot be modelled as a decision or counting problem and one has to come up with different measures to characterise the hardness of query evaluation. One approach that has been fruitfully applied to join evaluation is to study the worst-case output size as a measure of the hardness of a query. Atserias, Grohe and Marx [6] identified the fractional edge cover number of the join query as a crucial measure for lower bounding its worst-case output size. This bound was also shown to be optimal and is matched by so called “worst-case optimal” join evaluation algorithms, see [79, 95, 78, 63].

**Query Enumeration.** Another way of studying non-Boolean queries that is independent of the actual or worst-case output size is *query enumeration*. A query enumeration algorithm evaluates a non-Boolean query by reporting, one by one without repetition, the tuples in the query result. The crucial measure to characterise queries that are tractable w.r.t. enumeration is the delay between two output tuples. In the context of constraint satisfaction, the combined complexity, where the query as well as the database are given as input, has been considered. As the size of the query result might be exponential in the input size in this setting, queries that can be enumerated with *polynomial delay* and polynomial preprocessing are regarded as “tractable”. Classes of conjunctive queries that can be enumerated with polynomial delay have

## 1. Introduction

been identified in [24, 45]. However, a complete characterisation of conjunctive queries that are tractable in this sense is not in sight.

More relevant to the database setting, where one evaluates a small query against a large database, is the notion of *constant delay enumeration* introduced by Durand and Grandjean in [36]. The preprocessing time is supposed to be much smaller than the time needed to evaluate the query (usually, linear in the size of the database) and the delay between two output tuples may depend on the query, but not on the database. A lot of research has been devoted to this subject, where one usually tries to understand which structural restrictions on the query or on the database allow constant delay enumeration. For an introduction to this topic and an overview of the state-of-the-art the author refers the reader to the surveys [92, 91, 90].

Bagan, Durand and Grandjean [8] showed that acyclic conjunctive queries that are *free-connex* can be enumerated with constant delay after a linear time preprocessing phase (cf. [23] for a simplified proof of their result). They also showed that for self-join free acyclic conjunctive queries the free-connex property is essential by proving the following lower bound. Assume that multiplying two  $n \times n$  matrices cannot be done in time  $O(n^2)$ . Then the result of a self-join free acyclic conjunctive query that is not free-connex cannot be enumerated with constant delay after a linear time preprocessing phase.

It turns out that the notion of q-hierarchical conjunctive queries is a proper subclass of the free-connex conjunctive queries. Thus, there are queries that can be efficiently enumerated in the static setting but are hard to maintain under database updates.

Carmeli and Kröll generalized the notion of free-connex to UCQs [26] and they showed that they are tractable (in the sense that the result set can be enumerate with constant delay after a linear time preprocessing phase) and moreover, they showed that there are UCQs that are not free-connex but tractable. The problem of finding the subclass of UCQ, that are tractable in the static setting still remains open.

**First-order query evaluation on static databases.** In the thesis' second part, the query language FO+MOD, the extension of first-order logic FO with modulo-counting quantifiers of the form  $\exists^{i \bmod m} x \psi$ , expressing that the number of witnesses  $x$  that satisfy  $\psi$  is congruent to  $i$  modulo  $m$ , is studied.

Query evaluation algorithm is efficient if the update time is either constant or at most polylogarithmic ( $\log^c n$ ) in the size of the database. As a consequence, efficient query evaluation in the dynamic setting is only conceivable if the static problem (i.e., the setting without database updates) can be solved for Boolean queries in linear or pseudo-linear ( $n^{1+\varepsilon}$ ) time. Since this is not always possible, a short overview on known results about first-order query evaluation on static database is provided.

The problem of deciding whether a given database  $D$  satisfies a FO-sentence  $Q$  is AW[\*]-complete (parameterised by  $\|Q\|$ ) and it is therefore generally believed that the evaluation problem cannot be solved in time  $f(\|Q\|)\|D\|^c$  for any computable  $f$  and constant  $c$  (here,  $\|Q\|$  and  $\|D\|$  denote the size of the query and the database, respectively). For this reason, a long line of research focused on increasing classes of sparse instances ranging from databases of *bounded degree* [89] (where every domain element occurs only in a constant number of tuples in the database) to classes that

### 1.3. Recommended reading order

are *nowhere dense* [49]. In particular, Boolean first-order queries can be evaluated on classes of databases of bounded degree in linear time  $f(\|Q\|)\|D\|$ , where the constant factor  $f(\|Q\|)$  is 3-fold exponential in  $\|Q\|$  [89, 44], and Frick and Grohe [44] showed that the 3-fold exponential blow-up in terms of the query size is unavoidable assuming  $\text{FPT} \neq \text{AW}[*]$ .

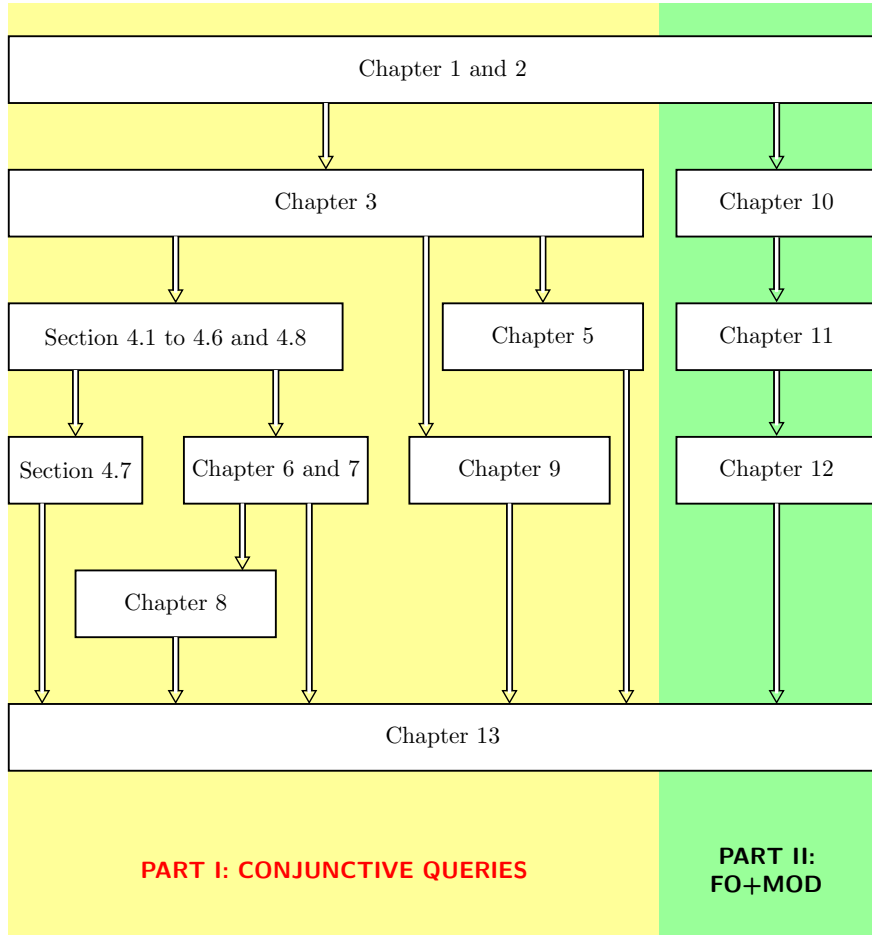
Durand and Grandjean [36] and Kazana and Segoufin [60] considered the task of enumerating the result of a  $k$ -ary first-order query on bounded degree databases and showed that after a linear time preprocessing phase the query result can be enumerated with constant delay. This result was later extended to classes of databases of bounded expansion [61]. Kazana and Segoufin [61] also showed that counting the number of result tuples of a  $k$ -ary first-order query on databases of bounded expansion (and hence also on databases of bounded degree) can be done in time  $f(\|Q\|)\|D\|$ . Segoufin and Vigny [93] proved an analogous result for classes of locally bounded expansion and pseudo-linear time  $f(\|Q\|)\|D\|^{1+\varepsilon}$ , and they also presented an algorithm for enumerating the query result with constant delay after pseudo-linear time preprocessing. These results were recently generalised to all nowhere dense classes of databases by Grohe and Schweikardt [50] and Schweikardt, Segoufin and Vigny [87]. Durand, Schweikardt and Segoufin [38] obtained analogous results for classes of databases of low degree (i.e., degree at most  $\|D\|^{o(1)}$ ).

### 1.3. Recommended reading order

Figure 1.1 shows a diagram that illustrates in which sequence the reader can read this thesis.

1. Introduction

Figure 1.1.: Recommended reading order.





## 2. Preliminaries

We write  $\mathbb{N}$  for the set of non-negative integers and let  $\mathbb{N}_{\geq 1} := \mathbb{N} \setminus \{0\}$  and  $[n] := \{1, \dots, n\}$  for all  $n \in \mathbb{N}_{\geq 1}$ . By  $2^S$  we denote the power set of a set  $S$ . For a partial function  $f$  we write  $\text{dom}(f)$  and  $\text{codom}(f)$  for the domain and the codomain of  $f$ , respectively. When writing  $\text{poly}(n)$ , we mean  $n^{O(1)}$  and when writing  $\exp(n)$  we mean  $2^{\text{poly}(n)}$ .

Let  $A$  be a set and let  $<$  be an order on  $A$  and let  $\bar{a}, \bar{b} \in A^k$  be two tuples of length  $k$  with  $\bar{a} = (a_1, \dots, a_k)$  and  $\bar{b} = (b_1, \dots, b_k)$ . A tuple  $\bar{a}$  is lexicographically smaller than  $\bar{b}$ , symb.  $\bar{a} <_{\text{lex}} \bar{b}$ , if there is a  $j \in [k]$  such that  $b_j < a_j$  and for all  $i < j$  we have  $b_i = a_i$ . We write  $\bar{a} \leq_{\text{lex}} \bar{b}$  if  $\bar{a} <_{\text{lex}} \bar{b}$  or  $\bar{a} = \bar{b}$ .

Let  $G$  be a graph. We will usually write  $V(G)$  to denote the vertex set of  $G$  and  $E(G)$  to denote the edge set of  $G$ . Let  $T$  be a tree.

**Principle of Inclusion and Exclusion.** In this thesis, one often has to compute the number of elements in the union of sets. To determine these cardinalities, we use the *principle of inclusion and exclusion* (see [25] for the statement and a proof). Let  $X$  be a set and  $(A_1, \dots, A_n)$  be a family of subsets of  $X$ . Then the number of elements of  $X$  which lie in none of the subsets  $A_i$  is

$$\sum_{\emptyset \neq \mathcal{I} \subseteq \{1, \dots, n\}} (-1)^{|\mathcal{I}|} \left| \bigcap_{i \in \mathcal{I}} A_i \right| + |X|. \quad (2.1)$$

Using (2.1) and setting  $X = A_1 \cup \dots \cup A_n$ , we obtain the following

$$\left| \bigcup_{i=1}^n A_i \right| = \sum_{\emptyset \neq \mathcal{I} \subseteq \{1, \dots, n\}} (-1)^{|\mathcal{I}|+1} \left| \bigcap_{i \in \mathcal{I}} A_i \right| \quad (2.2)$$

**Databases.** We fix a countably infinite set **dom**, the *domain* of potential database entries. Elements in **dom** are called *constants*. A *schema* is a finite set  $\sigma$  of relation symbols, where each  $R \in \sigma$  is equipped with a fixed *arity*  $\text{ar}(R) \in \mathbb{N}$ . Let us fix a schema  $\sigma = \{R_1, \dots, R_s\}$  and let  $r_i := \text{ar}(R_i)$  for  $i \in [s]$ . A *database*  $D$  of schema  $\sigma$  ( $\sigma$ -db for short) is of the form  $D = (R_1^D, \dots, R_s^D)$  where each  $R_i^D$  is a finite subset of  $\text{dom}^{r_i}$ . The *active domain*  $\text{adom}(D)$  of  $D$  is the smallest subset  $A$  of **dom** such that  $R_i^D \subseteq A^{r_i}$  for all  $i \in [s]$ .

**Updates.** We allow to update a given database of schema  $\sigma$  by inserting or deleting tuples as follows. An *insertion* command is of the form

$$\text{insert } R(a_1, \dots, a_r)$$

## 2. Preliminaries

for  $R \in \sigma$ ,  $r = \text{ar}(R)$  and  $a_1, \dots, a_r \in \mathbf{dom}$ . When applied to a  $\sigma$ -db  $D$ , the result is the updated  $\sigma$ -db  $D'$  with  $R^{D'} := R^D \cup \{(a_1, \dots, a_r)\}$  and  $S^{D'} := S^D$  for all  $S \in \sigma \setminus \{R\}$ . A *deletion* command is of the form

$$\text{delete } R(a_1, \dots, a_r)$$

for  $R \in \sigma$ ,  $r = \text{ar}(R)$  and  $a_1, \dots, a_r \in \mathbf{dom}$ . When applied to a  $\sigma$ -db  $D$ , the result is the updated  $\sigma$ -db  $D'$  with  $R^{D'} := R^D \setminus \{(a_1, \dots, a_r)\}$  and  $S^{D'} := S^D$  for all  $S \in \sigma \setminus \{R\}$ . Note that both types of commands may change the database's active domain.

**Sizes and Cardinalities.** The *cardinality*  $|D|$  of a  $\sigma$ -db  $D$  is defined as the number of tuples stored in  $D$ , i.e.,  $|D| := \sum_{R \in \sigma} |R^D|$ . The *size*  $\|D\|$  of  $D$  is defined as  $|\sigma| + |\text{adom}(D)| + \sum_{R \in \sigma} \text{ar}(R) \cdot |R^D|$  and corresponds to the size of a reasonable encoding of  $D$ . We will often write  $n$  to denote the cardinality  $|\text{adom}(D)|$  of  $D$ 's active domain.

**Dynamic Algorithms for Query Evaluation.** The algorithms take as input a  $k$ -ary query  $Q(x_1, \dots, x_k)$  and a  $\sigma$ -db  $D_0$ . For all query evaluation problems considered in this thesis, we aim at routines **preprocess**, **update** and **init** which achieve the following:

- upon input of  $Q(x_1, \dots, x_k)$  and  $D_0$ , **preprocess** builds a data structure  $\mathbf{D}$  which represents  $D_0$  and which is designed in such a way that it supports efficient evaluation of  $Q$  on  $D_0$ ,
- upon input of a command **update**  $R(a_1, \dots, a_r)$  (with **update**  $\in \{\text{insert}, \text{delete}\}$ ), calling **update** modifies the data structure  $\mathbf{D}$  such that it represents the updated database  $D$  and
- upon input of  $Q(x_1, \dots, x_k)$ , we denote with **init** the particular case of the routine **preprocess** upon input of the query  $Q(x_1, \dots, x_k)$  and the *empty* database  $D_\emptyset$ , where  $R^{D_\emptyset} = \emptyset$  for all  $R \in \sigma$ .

The *preprocessing time*  $t_p$  is the time used for performing **preprocess**; the *update time*  $t_u$  is the time used for performing an **update**; the *initialisation time*  $t_i$  is the time used for performing **init**. In all dynamic algorithms presented in this thesis, the **preprocess** routine for input of  $Q(x_1, \dots, x_k)$  and  $D_0$  will carry out the **init** routine for  $Q(x_1, \dots, x_k)$  and then perform a sequence of  $|D_0|$  update operations to insert all the tuples of  $D_0$  into the data structure. Consequently,  $t_p = t_i + |D_0| \cdot t_u$ . Thus, it is sufficient to consider the **update** and the **init** routine instead of a **preprocess** routine, if  $t_p = t_i + |D_0| \cdot t_u$ .

An enumeration routine **enumerate** is a routine that invokes an enumeration of the elements of a set  $M$  and afterwards outputs an end-of-enumeration message EOE. The delay of **enumerate** is defined as the maximum time used during a call of **enumerate**

- until the output of the first element in  $M$  (or the end-of-enumeration message EOE, if  $M = \emptyset$ ),
- between the output of two consecutive elements in  $M$  and

- between the output of the last element in  $M$  and the end-of-enumeration message EOE.

In the following,  $D$  will always denote the database that is currently represented by the data structure  $D$ . To solve the *enumeration problem under updates*, apart from the routines **preprocess** and **update**, we aim at a routine **enumerate** such that calling **enumerate** invokes an enumeration of all tuples, without repetition, that belong to the query result  $Q(D)$ . The *delay*  $t_d$  is the delay of **enumerate**.

To *test* if a given tuple belongs to the query result, we aim at a routine **test** which upon input of a tuple  $\bar{a} \in \text{dom}^k$  checks whether  $\bar{a} \in Q(D)$ . The *testing time*  $t_t$  is the time used for performing a **test**.

To solve the *counting problem under updates*, we aim at a routine **count** which outputs the cardinality  $|Q(D)|$  of the query result. The *counting time*  $t_c$  is the time used for performing a **count**.

To *answer* a *Boolean* query under updates, we aim at a routine **answer** that produces the answer **yes** or **no** of  $Q$  on  $D$ . The *answer time*  $t_{\text{ans}}$  is the time used for performing **answer**.

To *output the  $j$ th solution*, we aim at a routine **jth** which outputs upon input of a number  $j \in \mathbb{N}_{\geq 1}$  the  $j$ th tuple the **enumerate** routine would output if  $j$  is smaller or equal to the number of tuples in the result. The  *$j$ th time*  $t_j$  is the time used for performing **jth**.

To *enumerate the difference*, we aim at a routine **diff** which does the following. For a predefined  $\sigma$ -db  $D^-$  and a  $\sigma$ -db  $D$  where  $D$  is a database obtained from  $D^-$  after a couple of update steps, **diff** shall enumerate the sets  $Q(D^-) \setminus Q(D)$  and  $Q(D) \setminus Q(D^-)$ . The *diff-delay time*  $t_{di}$  is the delay for the routine **diff**.

To *prepare the learning of the result set*, we aim at a routine **learn** which prepares the algorithm for learning the query result as the training set in the polynomial regression model. The *learning time*  $t_l$  is the time for the routine **learn**.

Whenever speaking of a *dynamic algorithm*, we mean an algorithm that has routines **preprocess** and **update** and, depending on the problem at hand, at least one of the routines **enumerate**, **count**, **test**, **jth**, **diff**, **learn** and **answer**.

**Machine Model.** Following [32], we use Random Access Machines (RAMs) with  $O(\log n)$  word-size and a uniform cost measure to analyse our algorithms. We will assume that the RAM's memory is initialised to 0. In particular, if an algorithm uses an array, we will assume that all array entries are initialised to 0 and this initialisation comes at no cost (in real-world computers this can be achieved by using the *lazy array initialisation technique*, cf. e.g. [77]). A further assumption is that for every fixed dimension  $k \in \mathbb{N}_{\geq 1}$  we have an unbounded number of  $k$ -ary arrays  $A$  available such that for given  $(n_1, \dots, n_k) \in \mathbb{N}^k$  the entry  $A[n_1, \dots, n_k]$  at position  $(n_1, \dots, n_k)$  can be accessed in constant time.<sup>1</sup>

Throughout this thesis, we often adopt the view of *data complexity* and suppress factors that may depend on the query  $Q$  or the degree bound  $d$  (the degree bound

<sup>1</sup>While this can be accomplished in the RAM-model, for an implementation on real-world computers one would probably have to resort to replacing our use of arrays by using suitable designed hash functions.

## 2. Preliminaries

is only relevant for Part II) but not on the database  $D$ . E.g., “linear preprocessing time” means  $t_p \leq g(Q, d) \cdot \|D\|$  and “constant update time” means  $t_u \leq g(Q, d)$  for a function  $g$  with codomain  $\mathbb{N}$ .

**The yield command.** To keep algorithms simple and to simply analyse their running time we will often use the `yield` operator in our algorithms. This operator is also used in popular programming languages such as Python. The yield operator will be used in functions. As a consequence, these functions become iterable and an iteration over the functions works as follows. When we iterate over a function, the function starts in the first iteration and runs until it reaches the yield operator and then it outputs the value of the **yield** command. If the next value for the function is requested, the function continues from the point it stopped in the previous iteration and continues until it reaches the next yield and outputs the value of the yield operation. An example is given by the following Python code.

```
def f():
    for i in [1,2,3]:                // Iterate for i in {1,2,3}
        yield i*i
        print("Output for " + str(i))

for l in f():
    print(l)
    print("End")
```

With the yield operator, the function `f` is iterable. When we start with the for loop `l in f()` the function runs until `yield 1*1` is called and we receive 1 as output. Then we continue with the block of the for-loop `l in f()` and print `End` and then we continue with `f()` and print `Output for 1` and then stop at `yield 2*2` and so on. The for loop is finished when the execution of `f()` is finished, i.e., `Output for 3` will be printed at last. The output of the example is

```
1
End
Output for 1
4
End
Output for 2
9
End
Output for 3
```

The running time between two iteration steps of a function with yield is the maximum of the following running times:

- the time until we reach the first yield and
- the time between two yield operations and

- the time between a yield operation and the end of the function.

**The list data structure.** In the whole thesis, we consider the following dynamic data structure that represents a list which can be initialised and updated in constant time. For each list we have a type  $T$  which represents the set of possible items. We write  $\mathcal{L} = \langle b_1, \dots, b_\ell \rangle$  to denote a list of type  $T$  with  $\ell$  elements  $b_1, \dots, b_\ell \in T$ . For example,  $\langle 4, 6, 7 \rangle$  is a list of type  $\mathbb{N}$  with the elements 4, 6, 7. For a list  $\mathcal{L} = \langle b_1, \dots, b_\ell \rangle$  and an element  $a \in T$ , we write  $a \in \mathcal{L}$  if  $a \in \{b_1, \dots, b_\ell\}$  and  $a \notin \mathcal{L}$  if  $a \notin \{b_1, \dots, b_\ell\}$ .

We will first describe the data structure where an element  $a \in T$  can appear at most once in the list. Afterwards, we extend the data structure such that we can maintain a list where an element can appear multiple in the list.

The data structure for a list consists of an array  $A$ , a pointer *start* and a set of items which represents elements in the list. For every item  $i$  in the list we store the value  $i.a$  which represents the corresponding element of  $T$  and a pointer  $i.next$  that points to the successor element in the list, if it exists and to *nil* otherwise and a pointer  $i.prev$  that points to the previous element in the list, if it exists and to *nil* otherwise. The *next* and the *prev* pointer will help us to quickly insert and remove items in the list or *nil* if it is the empty list. The array  $A$  maps an element in  $T$  to the item in the list if it is part of the list and to *nil* otherwise. The pointer *start* points to the first item in the list.

For example, the data structure for the list  $\langle 4, 6, 7 \rangle$  (that is of type  $\mathbb{N}$ ) is the following. There are three items  $i_1, i_2, i_3$  with  $i_1.a = 4$  and  $i_2.a = 6$  and  $i_3.a = 7$  and  $i_j.next = i_{j+1}$  for all  $j \in \{1, 2\}$  and  $i_j.prev = i_{j-1}$  for all  $j \in \{2, 3\}$  and  $i_1.prev = i_3.next = \text{nil}$  and  $start = i_1$ . We have  $A[4] = i_1$ ,  $A[6] = i_2$ ,  $A[7] = i_3$  and  $A[n] = 0$  for  $n \in \mathbb{N} \setminus \{4, 6, 7\}$ .

To *initialise* the data structure for the empty list, we initialise the array  $A$  (note that by definition we can assume that  $A[a] = 0$  for all  $a \in T$  after the initialisation) and set the *start* pointer to *nil*. This can be done in  $O(1)$ .

The aim is for a list  $\mathcal{L}$  of type  $T$  to support two update operations  $\text{insert}(b, a)$  and  $\text{remove}(b)$  where  $b \in \mathcal{L} \cup \{\text{nil}\}$  and  $a \in T$ . On input of an update operation  $\text{insert}(b_k, a)$  with  $b_k \neq \text{nil}$  we insert  $a$  as the successor  $b_k$  to the list if  $a \notin \mathcal{L}$ , i.e., for the list  $\mathcal{L} = \langle b_1, \dots, b_\ell \rangle$  where  $k \leq \ell$ , the result after the update is  $\langle b_1, \dots, b_k, a, b_{k+1}, \dots, b_\ell \rangle$ . If  $b_k = \text{nil}$ , we insert  $a$  as the first element to the list, i.e., if the list  $\mathcal{L} = \langle b_1, \dots, b_\ell \rangle$  receives the update command  $\text{insert}(\text{nil}, a)$ , then the resulting list is  $\langle a, b_1, \dots, b_\ell \rangle$ . We will often use a variant of the insertion operation  $\text{insert}(a)$  where no previous element is specified. Then we will add the new item as the last element in the list, i.e.,  $\text{insert}(a)$  calls  $\text{insert}(b_\ell, a)$  where  $b_\ell$  is the last element of  $\mathcal{L}$  or  $b_\ell = \text{nil}$  if  $\mathcal{L}$  is empty. If  $\mathcal{L} = \langle b_1, \dots, b_\ell \rangle$  receives the update operation  $\text{remove}(b_k)$  for  $k \leq \ell$ , we remove the item with value  $b_k$  from the list, i.e., we obtain the list  $\langle b_1, \dots, b_{k-1}, b_{k+1}, \dots, b_\ell \rangle$ . The update operations will be implement as follows.

If the list  $\mathcal{L} = \langle b_1, \dots, b_\ell \rangle$  receives the update  $\text{insert}(b_k, a)$  and we implement the insertion algorithm as follows. If  $b_k \neq \text{nil}$ , we let  $i^-$  be the list item representing  $b_k$  and  $i^+$  be  $i^-.next$ . Note that  $i^+$  is *nil* if  $b_k$  is the last element in the list. We create a new list item  $i$  with  $i.a = a$ ,  $i.prev = i^-$  and  $i.next = i^+$ . We will also set  $i^-.next$  to  $i$  and, if  $i^+ \neq \text{nil}$ , the value  $i^+.prev$  to  $i$ . If we have the case that  $b_k = \text{nil}$ , we let  $i^+$  be the value of the *start* pointer and create a new item  $i$  with  $i.a = a$ ,  $i.prev = \text{nil}$

## 2. Preliminaries

and  $i.\text{next} = i^+$ . We will also set  $i^+.\text{prev}$  to  $i$  if  $i^+ \neq \text{nil}$ . Moreover, we set  $A[a]$  to the new item  $i$  in all cases. Note that we obtain a data structure for the list with the new element is inserted as above.

If the list  $\mathcal{L} = \langle b_1, \dots, b_\ell \rangle$  receives the update command  $\text{remove}(b_k)$ , let  $i$  be the item which represents  $b_k$  and let  $i^-$  be  $i.\text{prev}$  and  $i^+$  be  $i.\text{next}$ . If  $i^-$  is not  $\text{nil}$ , we set  $i^-. \text{next}$  to  $i^+$  and otherwise  $\text{start}$  to  $i^+$ . If  $i^+$  is not  $\text{nil}$ , we set  $i^+.\text{prev}$  to  $i^-$ . Furthermore, we set  $i.\text{prev} = i.\text{next} = \text{nil}$ , and we set  $A[a]$  to 0. Note that we obtain a data structure for the list where  $b_k$  is removed from the list as described above.

Since we have a constant number of “setting pointer to items” and lookup, insert and remove commands on  $A$  during both update commands, an update on the list takes time  $O(1)$ .

The data structure also supports a **lookup**, an **empty** and an **enumeration** routine. The **lookup** routine  $\text{lookup}(a)$  receives as input an element  $a \in T$  and outputs true if and only if  $a \in \mathcal{L}$ . To implement this operation, we simply output the result of testing  $A[a] \neq 0$ . This can be done in  $O(1)$ . The **empty** routine tests if the list is empty. This can be done in  $O(1)$  by testing if  $\text{start}$  is set to  $\text{nil}$ . The **enumerate** routine enumerates the elements in the list. This can be implemented as follows. We let  $i$  be the item that is pointed from  $\text{start}$  (if  $\text{start} = \text{nil}$ , we output **EOE** and terminate the enumerate operation). While  $i \neq \text{nil}$ , we yield  $i.a$  and set  $i$  to  $i.\text{next}$ . Afterwards, we output **EOE** and terminate the enumerate operation. The delay during the enumeration is  $O(1)$ .

To enrich that data structure for maintaining lists with multiple entries we use the following modifications. Instead of storing a single item in  $A[a]$ , we store in  $A[a]$  a pointer to a list  $\mathcal{L}_a$  of elements that consists of items in  $\mathcal{L}$  whose values are  $a$ . Note that since the pointer to an item in  $\mathcal{L}$  is unique, the elements in the list  $\mathcal{L}_a$  appears at most once and we can use the data structure described above for  $\mathcal{L}_a$ . Furthermore, we prune in the insertion procedure the condition that  $a \notin \mathcal{L}$  since an element can multiple appear in  $\mathcal{L}$ . When inserting an element  $a$  to the list  $\mathcal{L}$  we add the new item to the list  $\mathcal{L}_a$  (instead of setting  $A[a]$  to the item) and if we remove the item  $b_k$  from the list  $\mathcal{L}$  we remove the corresponding item from the list  $\mathcal{L}_{i.\text{value}}$  (instead of setting  $A[i.\text{value}] = 0$ ). To test if  $a \in \mathcal{L}$  for the **lookup** routine, we test if the list in  $A[a]$  is not empty and we do not have  $A[a] = 0$ . It is easy to verify that inserting and removing an element to/from  $\mathcal{L}$  and the **lookup** routine can be done in  $O(1)$ .

**AVL-trees for dynamic list representations.** In this thesis we will often use list representations for which we have to find the  $k$ th item in the list efficiently. To realise this, we will use the idea of Adelson-Velskii and Landis [3] that allows us to design the following data structure.

**Theorem 2.1** ([3],[65]). *There is a data structure that represents an arbitrary list of size  $N$  and allows the following operations in time  $O(\log(N))$ .*

1. *Upon input of a key, find the item in the list with the corresponding key,*
2. *upon input a number  $k$ , find the item in the  $k$ th position in the list,*
3. *insert an item (at a specific place) and*

4. delete a specified item.

We give an idea how the data structure in Theorem 2.1 works, see [65] for details. Additionally to a list we use a balanced binary tree. A balanced binary tree is a binary tree where for every node the following holds. The difference between the height of the left child and the height of the right node lies between  $-1$  and  $1$ . This gives a guarantee that the height of such a tree is in  $O(\log(N))$  where  $N$  is the size of the list.

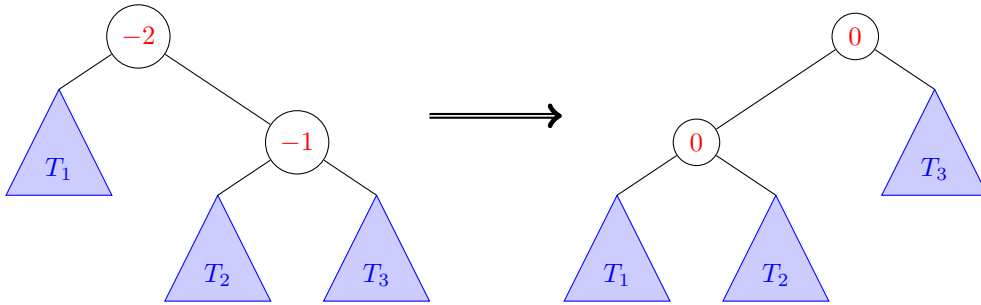
For every node, we store two pointer **LLINK** and **RLINK** that stores a pointer to the left and the right child, respectively. Additionally, we store for every node the fields:

- **KEY**: that represents a value of a list if the node is a leaf,
- **B**: that represents the balance factor of the node, i.e., the number of the height of the left child minus the number of the height of the right child,
- **RANK**: one plus the number of leafs in the left subtree.

The field **RANK** will help us to find the  $k$ th element in the list. In a balanced tree the field **B** for every node must be  $-1$ ,  $0$  or  $+1$ . After inserting or deleting an element to/from the list, the tree may not be balanced. To correct this, the algorithm makes some rotation on the tree (see Figure 2.1 and 2.2 for the rotations). In these figures, there are the two cases of an unbalanced trees shown. The red numbers show the balance factor of a node and the blue triangles depict that there is a subtree.

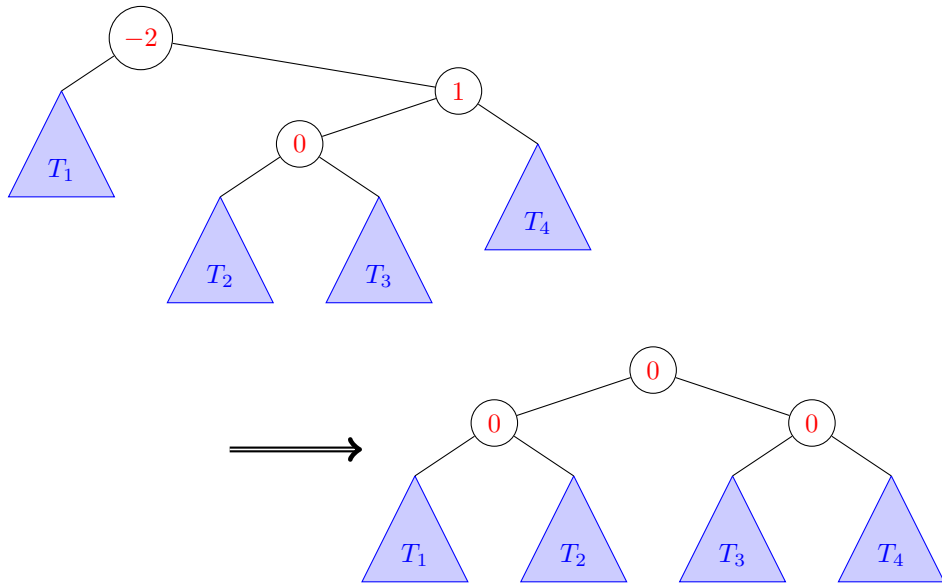
A detailed description of the mentioned algorithms is given in [65].

Figure 2.1.: Single-rotation.  $T_1$  and  $T_2$  are subtrees of height  $h$  and  $T_3$  is a subtree of height  $h + 1$ .



## 2. Preliminaries

Figure 2.2.: Double rotation. Each subtree  $T_1$ ,  $T_2$ ,  $T_3$  and  $T_4$  has height  $h$ .





**Part I.**

**Conjunctive Queries under  
Updates**



## 3. Main Theorems and Organisation of Part I

This chapter starts with a definition of syntax and semantics of conjunctive queries. Then the main results of this part are given and an organization of the first part.

### 3.1. An introduction to conjunctive queries

We fix a countably infinite set **var** of *variables*. An *atomic query* (for short: *atom*)  $\psi$  of schema  $\sigma$  is of the form  $Ru_1 \cdots u_r$  with  $R \in \sigma$ ,  $r = \text{ar}(R)$  and  $u_1, \dots, u_r \in \mathbf{var} \cup \mathbf{dom}$ . A *conjunctive query* (CQ, for short) of schema  $\sigma$  is of the form

$$\{(x_1, \dots, x_k, b_1, \dots, b_\ell) : \exists y_1 \cdots \exists y_\ell (\psi_1 \wedge \cdots \wedge \psi_d)\} \quad (3.1)$$

where  $k, \ell \in \mathbb{N}$ ,  $d \in \mathbb{N}_{\geq 1}$ ,  $\psi_j$  is an atomic query of schema  $\sigma$  for every  $j \in [d]$  and  $x_1, \dots, x_k, y_1, \dots, y_\ell$  are pairwise distinct elements in **var** and  $b_1, \dots, b_\ell \in \mathbf{dom}$ . A formula of the form  $\exists y_1 \cdots \exists y_\ell (\psi_1 \wedge \cdots \wedge \psi_d)$  is called conjunctive formula. *Join queries* are quantifier-free CQs, i.e., CQs of the form

$$\{(x_1, \dots, x_k, b_1, \dots, b_\ell) : (\psi_1 \wedge \cdots \wedge \psi_d)\}.$$

A CQ is called *self-join free* (or *non-repeating* or *simple*) if no relation symbol occurs more than once in the query. For a CQ  $Q$  of the form (3.1) we let  $\text{vars}(Q)$  ( $\text{cons}(Q)$ ) be the set of all variables (constants) occurring in  $Q$ , and we let  $\text{free}(Q) := \text{vars}(Q) \setminus \{y_1, \dots, y_\ell\} = \{x_1, \dots, x_k\}$  be the set of *free* variables. For an atom  $\psi$  let  $\text{vars}(\psi)$  ( $\text{cons}(\psi)$ ) be the set of all variables (constants) occurring in  $\psi$ . For every variable  $x \in \text{vars}(Q)$  we let  $\text{atoms}(x)$  be the set of all atoms  $\psi_j$  of  $Q$  such that  $x \in \text{vars}(\psi_j)$ .

The semantics of CQs are defined as usual: A *valuation* is a mapping  $\beta : \text{vars}(Q) \cup \mathbf{dom} \rightarrow \mathbf{dom}$  with  $\beta(a) = a$  for every  $a \in \mathbf{dom}$ . For every atom  $\psi = Ru_1 \cdots u_r$  in  $Q$  we write  $(D, \beta) \models \psi$  to denote that  $(\beta(u_1), \dots, \beta(u_r)) \in R^D$ . A valuation  $\beta$  is a *homomorphism* from  $Q$  to a  $\sigma$ -db  $D$  if for every atom  $\psi$  in  $Q$  we have  $(D, \beta) \models \psi$ . We sometimes write  $(D, \beta) \models Q$  to indicate that  $\beta$  is a homomorphism from  $Q$  to  $D$ . The *query result*  $Q(D)$  of a  $k$ -ary CQ  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$  on a  $\sigma$ -db  $D$  is defined as the set  $\{(\beta(x_1), \dots, \beta(x_k), b_1, \dots, b_\ell) : \beta \text{ is a homomorphism from } Q \text{ to } D\}$ . Clearly,  $Q(D) \subseteq \text{adom}(D)^{k+\ell}$ . For a tuple  $\bar{a} \in \mathbf{dom}^k$  we sometimes write  $D \models Q[\bar{a}]$  to indicate that there is a homomorphism  $\beta : Q \rightarrow D$  with  $\bar{a} = (\beta(x_1), \dots, \beta(x_k), b_1, \dots, b_\ell)$ . We sometimes write  $(D, \alpha) \models \varphi$  for a conjunctive formula  $\varphi$  to indicate that  $(D, \beta) \models Q$  for  $Q := \{(x_1, \dots, x_k, b_1, \dots, b_\ell) : \varphi\}$  where  $\{x_1, \dots, x_k\}$  is the set of free variables

### 3. Main Theorems and Organisation of Part I

in  $\varphi$ , i.e., the set of variables occurring in  $\varphi$  that are not quantified and  $\{b_1, \dots, b_\ell\}$  the set of constants that appear in  $\varphi$ .

A *Boolean* CQ is a CQ  $Q$  with  $\text{free}(Q) = \emptyset$ . As usual, for Boolean CQs  $Q$  we will write  $Q(D) = \text{yes}$  instead of  $Q(D) \neq \emptyset$ , and  $Q(D) = \text{no}$  instead of  $Q(D) = \emptyset$ .

Two  $k$ -ary queries  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$  and  $Q'(x'_1, \dots, x'_k, b_1, \dots, b_\ell)$  are *equivalent* ( $Q \equiv Q'$ , for short) if  $Q(D) = Q'(D)$  for every  $\sigma$ -db  $D$ .

The *size*  $\|Q\|$  of a CQ  $Q$  is defined as the length of  $Q$  when viewed as a word over the alphabet  $\sigma \cup \mathbf{var} \cup \mathbf{dom} \cup \{\exists, \wedge, (, ), \{, \}\}$ . For a  $k$ -ary CQ  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$  and a  $\sigma$ -db  $D$ , the *cardinality of the query result* is the number  $|Q(D)|$  of tuples in  $Q(D)$ .

The *arity* of a CQ  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$  is the number  $k + \ell$ .

We write  $\text{poly}(Q)$  as an abbreviation for  $\text{poly}(\|Q\|)$  and we write  $\text{exp}(Q)$  as an abbreviation for  $\text{exp}(\|Q\|)$ .

**UCQs.** A  $k$ -ary *union of conjunctive queries* ( $k$ -ary UCQ) is of the form  $Q_1(\bar{u}_1) \cup \dots \cup Q_d(\bar{u}_d)$  where  $d > 1$  and  $Q_i(\bar{u}_i)$  is a  $k$ -ary CQ of schema  $\sigma$  for every  $i \in [d]$ . The query result of such a  $k$ -ary UCQ  $Q$  on a  $\sigma$ -db  $D$  is  $Q(D) := \bigcup_{i=1}^d Q_i(D)$ . For a  $k$ -ary query  $Q$  we write  $\text{vars}(Q)$  (and  $\text{cons}(Q)$ ) to denote the set of all variables (and constants) that occur in  $Q$ . Clearly,  $Q(D) \subseteq (\text{adom}(Q) \cup \text{cons}(Q))^k$ .

## 3.2. Main theorems and organisation

The aim of the thesis' first part is to investigate the following problem. For which class of conjunctive queries, we can construct a data structure that represents the result set on a database that is tractable under updates and for which we can design fast enumeration, answering, testing, difference,  $j$ th routines. It turns out that the class of conjunctive queries for which we can realise this is the class of q-hierarchical conjunctive queries, which are defined as follows.

**Definition 3.1.** A CQ  $Q$  is q-hierarchical if for any two variables  $x, y \in \text{vars}(Q)$  the following is satisfied:

- (i)  $\text{atoms}(x) \subseteq \text{atoms}(y)$  or  $\text{atoms}(x) \supseteq \text{atoms}(y)$  or  $\text{atoms}(x) \cap \text{atoms}(y) = \emptyset$ ,  
and
- (ii) if  $\text{atoms}(x) \subsetneq \text{atoms}(y)$  and  $x \in \text{free}(Q)$ , then  $y \in \text{free}(Q)$ .

Minimal examples for queries that are not q-hierarchical are

$$Q_{S-E-T} := \{(x, y, z) : (Sx \wedge Exy \wedge Ty)\} \quad (3.2)$$

and

$$Q_{E-T} := \{(x) : \exists y (Exy \wedge Ty)\} \quad (3.3)$$

because they do not satisfy condition (i) and (ii), respectively. Note that some variations of  $Q_{E-T}$  such as the join query  $\{(x, y) : (Exy \wedge Ty)\}$ , the Boolean query  $\{() : \exists x \exists y (Exy \wedge Ty)\}$  and the query  $\{(y) : \exists x (Exy \wedge Ty)\}$  are q-hierarchical.

### 3.2. Main theorems and organisation

The notion of *q-hierarchical* conjunctive queries is related to the *hierarchical* property that has already played a central role for efficient query evaluation in various contexts. It has been introduced by Dalvi and Suciu in [34] to characterise the Boolean CQs that can be answered in polynomial time on probabilistic databases. They obtained a dichotomy stating for self-join free queries that the complexity of query evaluation on probabilistic databases is in PTIME for hierarchical queries and #P-complete for non-hierarchical queries. Fink and Olteanu [42] generalised the notion and the dichotomy result to non-Boolean queries and to queries using negation. In the different context of query evaluation on massively parallel architectures, Koutris and Suciu [69] considered hierarchical join queries and singled out a subclass of so-called tall-flat queries as exactly those queries that can be computed with only one broadcast step in their *Massively Parallel* model of query evaluation. For further information on the various uses of the hierarchical property the author refers the reader to [42].

The definition of hierarchical queries relies on the following notion. Consider a CQ  $Q$  of the form (3.1). Dalvi and Suciu [34] call a Boolean CQ  $Q$  *hierarchical* iff the condition

$$(*) : \text{atoms}(x) \subseteq \text{atoms}(y) \text{ or } \text{atoms}(x) \supseteq \text{atoms}(y) \text{ or } \text{atoms}(x) \cap \text{atoms}(y) = \emptyset$$

is satisfied by all variables  $x, y \in \text{vars}(Q)$ . An example for a hierarchical Boolean CQ is  $\{() : \exists x \exists y \exists z \exists y' \exists z' (Rxyz \wedge Rxyz' \wedge Exy \wedge Exy')\}$ .

In [69], Koutris and Suciu transferred the notion to join queries  $Q$ , which they call hierarchical iff condition  $(*)$  is satisfied by all variables  $x, y \in \text{vars}(Q)$ . In [42], Fink and Olteanu introduced a slightly different notion for a more general class of queries. Translated into the setting of CQs, their notion (only) requires that condition  $(*)$  is satisfied by all *quantified* variables, i.e., variables  $x, y \in \text{vars}(Q) \setminus \text{free}(Q)$ . Obviously, both notions coincide on *Boolean* CQs, but on *join queries* Koutris and Suciu's notion is more restrictive than Fink and Olteanu's notion (according to which *all* quantifier-free CQs are hierarchical). For example, the join query  $Q_{S-E-T}$  from Equation (3.2) is hierarchical w.r.t. Fink and Olteanu's notion, and non-hierarchical w.r.t. Koutris and Suciu's notion.

To ensure tractability of a conjunctive query in that setting, we will require that its quantifier-free part is hierarchical in Koutris and Suciu's notion and, additionally, the quantifiers respect the query's hierarchical form. Such queries are the *q-hierarchical* queries.

Note that a *Boolean* CQ is q-hierarchical if and only if it is hierarchical, and a join query is q-hierarchical if and only if it is hierarchical w.r.t. Koutris and Suciu's notion.

This notion of q-hierarchical queries is also related to the notion of factorized databases from [82]. Every query result of a q-hierarchical query can be identified as a factorized database but not every factorized database can be described by a q-hierarchical query. Moreover, factorized database have only been considered in the static setting, i.e., in the setting without updates.

In the first part of this thesis we will prove the following theorems for q-hierarchical CQs:

### 3. Main Theorems and Organisation of Part I

**Theorem 3.2.** *There is a dynamic algorithm that receives a Boolean q-hierarchical conjunctive query  $Q$  and a  $\sigma$ -db  $D_0$ , and computes within  $t_i = \text{poly}(Q)$  initialisation time and  $t_p = \text{poly}(Q) \cdot O(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)$  and allows to answer  $Q(D)$  in time  $t_{ans} = O(1)$  where  $D$  is the current database.*

For non-Boolean q-hierarchical CQs we receive the following theorem.

**Theorem 3.3.** *There is a dynamic algorithm that receives a non-Boolean q-hierarchical conjunctive query  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$  and a  $\sigma$ -db  $D_0$  and computes within  $t_i = \text{poly}(Q)$  initialisation time and  $t_p = \text{poly}(Q) \cdot O(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)$  and allows to*

- (a) *test for an input tuple  $\bar{a} \in \text{dom}^{k+\ell}$  if  $\bar{a} \in Q(D)$  within time  $t_t = \text{poly}(Q)$ , and*
- (b) *enumerate  $Q(D)$  with delay  $t_d = \text{poly}(Q)$ , and*
- (c) *enumerate the tuples that join and omit the result set  $Q(D)$  after an arbitrary number of updates on a previous chosen state of the database received with delay  $t_{di} = \text{poly}(Q)$  and*
- (d) *compute the cardinality  $|Q(D)|$  in time  $t_c = O(1)$*

where  $D$  is the current database.

The result from Theorem 3.2 and Theorem 3.3 (a) (b) and (d) are results from [16]. The proof of these theorems are different to the proofs in [16]. The result of Theorem 3.3 (c) is a new result.

A proof for Theorem 3.2 and Theorem 3.3 is given in Chapter 4 (except for part (d) of Theorem 3.3. This is proven in Section 6.5). The data structure used in these theorems are the same and it is described in Section 4.1. In Section 4.3 a proof for Theorem 3.2 is given. Then, we show how to use the data structure for the testing problem (see Section 4.4) and for the enumeration problem (see Section 4.6) and for the difference problem (see Section 4.7). It turns out that there are non q-hierarchical queries for which testing can be done efficiently under updates. This superclass of the q-hierarchical CQs are called t-hierarchical queries and are defined as follows.

**Definition 3.4.** *A CQ  $Q$  is t-hierarchical if the following is satisfied.*

- (i) *for all  $x, y \in \text{vars}(Q) \setminus \text{free}(Q)$ , we have*  

$$\text{atoms}(x) \subseteq \text{atoms}(y) \text{ or } \text{atoms}(y) \subseteq \text{atoms}(x) \text{ or } \text{atoms}(x) \cap \text{atoms}(y) = \emptyset,$$
*and*
- (ii) *for all  $x \in \text{free}(Q)$  and all  $y \in \text{vars}(Q) \setminus \text{free}(Q)$ , we have*  

$$\text{atoms}(x) \cap \text{atoms}(y) = \emptyset \text{ or } \text{atoms}(y) \subseteq \text{atoms}(x).$$

### 3.2. Main theorems and organisation

The queries  $Q_{E-E-R} := \{(x, y) : \exists v_1 \exists v_2 \exists v_3 (Exv_1 \wedge Eyv_2 \wedge Rxyv_3)\}$  and  $Q_{S-E-T}$  from Equation (3.2) are examples of queries that are t-hierarchical but not q-hierarchical. In fact, the only difference between q-hierarchical CQs and t-hierarchical CQs relies on condition (i) in their definitions. While

$$\text{atoms}(x) \subseteq \text{atoms}(y) \text{ or } \text{atoms}(x) \supseteq \text{atoms}(y) \text{ or } \text{atoms}(x) \cap \text{atoms}(y) = \emptyset$$

holds for all  $x, y \in \text{vars}(Q)$  for q-hierarchical queries, the fact has to hold only for  $x, y \in \text{vars}(Q) \setminus \text{free}(Q)$ . The second condition in Definition 3.1 (for q-hierarchical queries) is equivalent to the second condition in Definition 3.4 which follows from the following lemma.

**Lemma 3.5.** *For all CQs  $Q$  for which condition (i) from Definition 3.1 holds are the following two conditions equivalent:*

- (a) *For all  $x, y \in \text{vars}(Q)$  the following holds: if  $\text{atoms}(x) \subsetneq \text{atoms}(y)$  and  $x \in \text{free}(Q)$ , then  $y \in \text{free}(Q)$ .*
- (b) *For all  $x \in \text{free}(Q)$  and all  $y \in \text{vars}(Q) \setminus \text{free}(Q)$ , we have  $\text{atoms}(x) \cap \text{atoms}(y) = \emptyset$  or  $\text{atoms}(y) \subseteq \text{atoms}(x)$ .*

*Proof.* Condition (a) holds, if and only if

$$\begin{aligned} &\text{for all } x, y \in \text{vars}(Q) \text{ is } x \notin \text{free}(Q) \text{ or } y \in \text{free}(Q) \\ &\text{or it holds not } \text{atoms}(x) \subsetneq \text{atoms}(y). \end{aligned} \quad (3.4)$$

From condition (i) of Definition 3.1 follows that it holds not  $\text{atoms}(x) \subsetneq \text{atoms}(y)$ , if and only if  $\text{atoms}(x) \supseteq \text{atoms}(y)$  or  $\text{atoms}(x) \cap \text{atoms}(y) = \emptyset$ . Thus, (3.4) is equivalent to

$$\begin{aligned} &\text{for all } x, y \in \text{vars}(Q) \text{ is } x \notin \text{free}(Q) \text{ or } y \in \text{free}(Q) \\ &\text{or it holds } \text{atoms}(x) \supseteq \text{atoms}(y) \text{ or } \text{atoms}(x) \cap \text{atoms}(y) = \emptyset \end{aligned} \quad (3.5)$$

It is straightforward to see that (3.5) and (b) are equivalent.  $\square$

We obtain the following theorem for the testing problem:

**Theorem 3.6** ([21]). *There is a dynamic algorithm that receives a t-hierarchical k-ary CQ  $Q$  and a  $\sigma$ -db  $D_0$ , and computes within  $t_p = \text{poly}(Q) \cdot O(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)$  and allows to test for an input tuple  $\bar{a} \in \text{dom}^k$  if  $\bar{a} \in Q(D)$  within time  $t_t = \text{poly}(Q)$ , where  $D$  is the current database.*

A proof for Theorem 3.6 is given in Chapter 5. Theorem 3.6 has already been published in [21] and the proof is based on the proof in [21].

In Chapter 6, we enrich q-hierarchical conjunctive queries to queries with aggregates.

In practice, aggregates can be used in every commercial implementation of SQL. Furthermore there are a lot of applications for queries with aggregates such as mobile computing [12], global information systems [71], stream data analysis [35], constraint databases [14]. That is the reason why lot of research has been done on evaluating

### 3. Main Theorems and Organisation of Part I

queries with aggregates [27, 52, 94, 31]. Moreover, aggregates can also be used to extend logics [56, 13, 40]. Only little is known about the enumeration problem for queries with aggregates in the static setting. In [10] it is shown how to enumerate the query result where the queries are factorized queries with aggregates that are very close to the q-hierarchical queries. In [10] the static setting has been considered.

In Chapter 6 we allow to operate over the result set with aggregates. In Section 6.1 there are examples for queries with aggregates and in Section 6.3 a formal definition of the syntax and semantics of such queries is given. It turns out that queries with aggregates are convenient to solve the counting problem for q-hierarchical conjunctive queries without aggregates, i.e., we can use queries with aggregates to prove Theorem 3.3(d) (see Section 6.5). In Section 6.6 and Section 6.7 we show how we can if the result of Theorem 3.3 from q-hierarchical conjunctive queries without aggregates to queries with aggregates and we obtain a proof for the following theorem:

**Theorem 3.7.** *There is a dynamic algorithm that receives a q-hierarchical CQ with aggregates  $Q$  with aggregation time  $t_a$  and a  $\sigma$ -db  $D_0$ , and computes within  $t_i = O(1)$  initialisation time and  $t_p = \text{poly}(Q)t_a O(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)t_a$  and allows to*

- (a) *enumerate  $Q(D)$  with delay  $t_d = \text{poly}(Q)$ ,*
- (b) *test for an input tuple  $\bar{a} \in \text{dom}^k$  if  $\bar{a} \in Q(D)$  within time  $t_t = \text{poly}(Q)$ ,*
- (c) *enumerate the tuples that join and omit the result set  $Q(D)$  after an arbitrary number of updates on a previous chosen state the database received with delay  $t_{di} = \text{poly}(Q)$  and*
- (d) *compute the cardinality  $|Q(D)|$  in time  $t_c = O(1)$*

*where  $D$  is the current database.*

This result is yet unpublished and a new result.

The aggregation time  $t_a$  of a query is the time it takes to update the values computed by the aggregation functions during an update step.

In Chapter 7 we show how to use the query result of a database query to prepare learning a polynomial function where the query result is taken as training data.

In the last years, the research interest for connecting database with machine learning has been growing [1].

Most machine learning systems and database systems do not support a connection between a database and the machine learning algorithm. To learn with data from a query result, one has to store the query result from a database system in a file that can be used as input for a machine learning algorithm. Such a data exchange is very expensive due to the repetition of data blocks e.g. in the result of join queries. That is the reason why it might be convenient to do some computation step of the machine learning algorithm in the database system, the so-called in-database learning. In-database learning was considered by Schleich and Olteanu [86] for factorized databases and Khamis et.al. [62] introduced a unified framework for in-database learning of



### 3.2. Main theorems and organisation

a class of statistical learning models in relation databases. In this thesis there is a in-database learning solution for regression models in relational databases on q-hierarchical queries under updates. Every time a tuple is inserted into or deleted from the database, we can appropriately modify the data structure of the in-database learning system.

Further work has been done in the task of Machine Learning and databases [64, 58, 84, 5].

We obtain the following new, unpublished, theorem.

**Theorem 3.8.** *For every  $d \in \mathbb{N}_{\geq 1}$  there is a dynamic algorithm that receives a  $q$ -hierarchical conjunctive query  $Q$  and a  $\sigma$ -db  $D_0$  and computes within  $t_i = O(\text{poly}(Q)^{2d+1})$  initialisation time and  $t_p = \text{poly}(Q)^{2d+1} \cdot O(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)^{2d+1}$  and allows to prepare the learning of the parameters of a polynomial hypothesis function  $h_\Theta$  of degree  $d$  in time  $t_l = O(\text{poly}(Q)^{2d+1})$  where  $D$  is the current database.*

Afterwards we consider the question if there are more routines we can do, if we allow logarithmic update time. When allowing logarithmic update time, one can organize the data structure from Theorem 3.3, such that we can enumerate the tuples in lexicographical order. This means, if a tuple  $\bar{a}$  will be enumerated before  $\bar{b}$ , then  $\bar{a} <_{\text{lex}} \bar{b}$ . Furthermore, one can upon input of a natural number output the  $j$ th tuple the enumeration routine would output and, upon input of a tuple in the result set, receive a natural number  $j$  such that the tuple is the  $j$ th tuple in the enumeration. This is stated in the next Theorem.

**Theorem 3.9.** *There is a dynamic algorithm that receives a  $q$ -hierarchical  $k$ -ary CQ with aggregates  $Q$  with aggregation time  $t_a$  and a  $\sigma$ -db  $D_0$  of size  $\|D_0\|$ , and computes within preprocessing time  $\text{poly}(Q) \cdot O(\|D_0\| \log(\|D_0\|)) t_a$  a data structure that can be updated in time  $\text{poly}(Q) \cdot O(\log(\|D\|)) t_a$  and allows the following:*

- (a) output the query result size  $|Q(D)|$  in time  $O(1)$  and
- (b) enumerate the tuples in  $Q(D)$  in lexicographical order with delay  $\text{poly}(Q)$  and
- (c) test for an input tuple  $\bar{a} \in \text{dom}^{k+\ell}$  if  $\bar{a} \in Q(D)$  within time  $t_t = \text{poly}(Q)$  and
- (d) upon input of a tuple  $\bar{b} \in \text{dom}^{k+\ell}$  output the tuple

$$\max \{ \bar{a} \in Q(D) : \bar{a} \leq \bar{b} \}$$

if it exists, or a *SmallerThanMinimum*-message otherwise in time  $t_l = \text{poly}(Q)$  and

- (e) upon input of an arbitrary number  $j$ , take time  $\text{poly}(Q) \cdot O(\log(\|D\|))$  to immediately output the  $j$ th tuple that the enumeration procedure of (b) would output and
- (f) upon input of an arbitrary tuple  $\bar{a} \in Q(D)$ , take time  $\text{poly}(Q) \cdot O(\log(\|D\|))$  to immediately output the number  $j$  such that  $\bar{a}$  is the  $j$ th tuple the enumeration procedure of (b) would output,

### 3. Main Theorems and Organisation of Part I

where  $D$  is the current database.

The result of this theorem is new and yet unpublished. To solve the routines in Part (a), (b), (c) and (d) we use a variant of the data structure from Theorem 3.3. The variant is described in Remark 4.9. Using this variant, we obtain that the query result will be enumerated in lexicographical order, when the enumeration algorithm from Theorem 3.3 is used (see Remark 4.15). Furthermore, we obtain the routine in Part (d) which is shown in Section 4.8.

For the routines in Part (e) and in Part (f) we enrich the data structure from Theorem 3.3 by another way. Note that, in particular, it is not necessary that the enumeration algorithm enumerates the tuples in lexicographical order for solving the  $j$ th problem and the  $j$ th-reverse problem. In Chapter 8, it is described how to enrich the data structure for the  $j$ th problem and the  $j$ th reverse problem and how these routines work.

To output the  $j$ th solution of the query result allows us to have random access on the query result by using any generation of integers. For the static setting Bagan, Durand, Grandjean and Olive [9] proved that for a FO query computing a random tuple of the query result can be done in average constant time after a linear preprocessing phase. Furthermore, Bagan [7] proved that, for a monadic second order query over classes of structures with bounded treewidth, the  $j$ th tuple in the query result with respect to an order  $<$  can be computed in time  $\log \|D\| + j$  after a linear time preprocessing phase.

In Chapter 9 we study unions of conjunctive queries.

We receive Theorem 3.10 for testing and Theorem 3.11 for enumerating.

**Theorem 3.10** ([21]). *There is a dynamic algorithm that receives a  $t$ -hierarchical  $k$ -ary UCQ  $Q$  and a  $\sigma$ -db  $D_0$ , and computes within  $t_p = \text{poly}(Q) \cdot O(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)$  and allows to test for an input tuple  $\bar{a} \in \text{dom}^k$  if  $\bar{a} \in Q(D)$  within time  $t_t = \text{poly}(Q)$ . Furthermore, the algorithm allows to answer a  $t$ -hierarchical Boolean UCQ within time  $t_{ans} = O(1)$ .*

**Theorem 3.11** ([21]). *There is a dynamic algorithm that receives a  $q$ -hierarchical  $k$ -ary UCQ  $Q$  and a  $\sigma$ -db  $D_0$ , and computes within  $t_p = \text{poly}(Q) \cdot O(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)$  and allows to enumerate  $Q(D)$  with delay  $t_d = \text{poly}(Q)$ .*

A proof for Theorem 3.10 and Theorem 3.11 is given in Section 9.1. The results of these theorems has already been published in [21]. The proof of Theorem 3.10 is closely based on the proof in [21]. For Theorem 3.11 an alternative proof is given in Section 9.1. In Section 9.2 we look at the  $j$ th problem for unions of conjunctive queries. We identify a class of unions of conjunctive queries for which we obtain a result, the strongly exhaustively  $q$ -hierarchical UCQs. A UCQ  $Q_1 \cup \dots \cup Q_m$  is strongly exhaustively  $q$ -hierarchical if  $\bigcap_{i \in [m]} Q_i$  is  $q$ -hierarchical. The corresponding theorem is the following.

### 3.2. Main theorems and organisation

**Theorem 3.12.** *There is a dynamic algorithm that receives a strongly exhaustively  $q$ -hierarchical  $k$ -ary UCQ  $Q$  and a  $\sigma$ -db  $D_0$  of size  $\|D_0\|$ , and computes within preprocessing time  $\exp(Q) \cdot O(\|D_0\| \log(\|D_0\|))$  a data structure that can be updated in time  $\exp(Q) \cdot O(\log(\|D\|))$  and allows the following:*

- (a) *enumerate the tuples in  $Q(D)$  with delay  $\text{poly}(Q)$  and*
- (b) *upon input of an arbitrary number  $j$ , take time  $\exp(Q) \cdot O(\log(\|D\|))$  to immediately output the  $j$ th tuple that the enumeration procedure of (a) would output,*

*where  $D$  is the current database.*

It turns out that if a conjunctive query is not  $q$ -hierarchical, we cannot maintain the query result under updates. This thesis mainly considers the upper bound, the lower bounds can be found in [16]. These bounds are conditioned on the OMv-conjecture, a conjecture on the hardness of online matrix-vector multiplication that characterises the hardness of dynamic problems [57]. In [16], we obtain the following dichotomies, which are stated from the perspective of data complexity (i.e., the query is regarded to be fixed) and hold for any fixed  $\varepsilon > 0$  under the OMv-conjecture. By  $n$  we denote in the remainder of this section the size of the active domain of the current database  $D$ . For the enumeration problem we restrict our attention to self-join free CQs, where every relation occurs only once in the query.

**Theorem 3.13** ([16]). *Fix a number  $\varepsilon > 0$  and a self-join free conjunctive query  $Q$ . If  $Q$  is not  $q$ -hierarchical, there is no algorithm with arbitrary preprocessing time and  $O(n^{1-\varepsilon})$  update time that enumerates  $Q(D)$  with  $O(n^{1-\varepsilon})$  delay, unless the OMv-conjecture fails.*

For *Boolean* CQs we obtain a lower bound for *all* queries, i.e., also for queries that are not self-join free. To state the result, we need the standard notion of a homomorphic core. A *homomorphism* from a CQ  $Q(x_1, \dots, x_k)$  to a CQ  $Q(y_1, \dots, y_k)$  is a mapping  $h$  from  $\text{vars}(Q)$  to  $\text{vars}(Q)$  such that  $h(x_i) = y_i$  for all  $i \in [k]$  and if  $Ru_1 \cdots u_r$  is an atom of  $Q$ , then  $Rh(u_1) \cdots h(u_r)$  is an atom of  $Q$ . The *homomorphic core* (for short, *core*) of a conjunctive query  $Q$  is a minimal subquery  $Q'$  of  $Q$  such that there is a homomorphism from  $Q$  to  $Q'$ , but no homomorphism from  $Q'$  to a proper subquery of  $Q$ . By Chandra and Merlin's homomorphism theorem, every CQ  $Q$  has a unique (up to isomorphism) core  $Q'$  and  $Q'(D) = Q(D)$  for all databases  $D$  (cf., e.g., [2]). While self-join free queries are their own cores, the situation is different for general CQs. Consider, for example, the queries

$$\begin{aligned} Q_1 &:= \{() : \exists x \exists y (Exx \wedge Exy \wedge Eyy)\} \quad \text{and} \\ Q_2 &:= \{() : \exists x (Exx)\}. \end{aligned}$$

Here,  $Q_2$  is a core of  $Q_1$  and thus  $Q_1(D) = Q_2(D)$  for every database  $D$ . However,  $Q_2$  is  $q$ -hierarchical, whereas  $Q_1$  is not. The next lower bound theorem states that the result of a Boolean conjunctive query cannot be maintained efficiently if the query's core is not  $q$ -hierarchical.

### 3. Main Theorems and Organisation of Part I

**Theorem 3.14** ([16]). *Fix a number  $\varepsilon > 0$  and a Boolean conjunctive query  $Q$ . If the homomorphic core of  $Q$  is not  $q$ -hierarchical, then there is no algorithm with arbitrary preprocessing time and  $t_u = O(n^{1-\varepsilon})$  update time that answers  $Q(D)$  in time  $t_{ans} = O(n^{2-\varepsilon})$ , unless the OMv-conjecture fails.*

Let us now turn to the problem of computing the cardinality  $|Q(D)|$  of the result of a query  $Q = \{(x_1, \dots, x_k) : \varphi\}$ . From the Theorems 3.3 and 3.14 we know that we can efficiently decide whether  $|Q(D)| > 0$  if and only if the homomorphic core of  $\{() : \exists x_1 \dots \exists x_k \varphi\}$  is  $q$ -hierarchical. The complexity of actually counting the number of tuples in  $Q(D)$ , however, depends on whether the core of the query  $Q(x_1, \dots, x_k)$  itself (rather than the core of its Boolean version  $\{() : \exists x_1 \dots \exists x_k \varphi\}$ ) is  $q$ -hierarchical. As in the Boolean case, the next theorem (together with Theorem 3.3) implies a dichotomy for all conjunctive queries. One difference is that we have to additionally rely on the OV-conjecture.

**Theorem 3.15** ([16]). *Fix a number  $\varepsilon > 0$  and a conjunctive query  $Q$ . If the homomorphic core of  $Q$  is not  $q$ -hierarchical, then there is no algorithm with arbitrary preprocessing time and  $t_u = O(n^{1-\varepsilon})$  update time that computes  $|Q(D)|$  in time  $t_c = O(n^{1-\varepsilon})$ , assuming the OMv-conjecture and the OV-conjecture.*

Last but not least, there is a new lower bound shown in this thesis for the  $j$ th problem, which can be shown by a reduction to the previous lower bound theorem:

**Theorem 3.16.** *Fix a number  $\varepsilon > 0$  and a self-join free conjunctive query  $Q$ . If  $Q$  is not  $q$ -hierarchical, then there is no algorithm with arbitrary preprocessing time and  $O(n^{1-\varepsilon})$  update time that solves the  $j$ th problem in time  $O(n^{1-\varepsilon})$ , unless the OMv-conjecture fails.*

A proof is given in Section 8.5.

## 4. Answering q-hierarchical conjunctive queries under updates

This chapter is devoted to the proof of Theorem 3.2 and Theorem 3.3, except for part (d) in Theorem 3.3. In Section 4.1 there is a definition of the data structure that will be used in both theorems. In Section 4.2 it is described how to initialise and update the data structure. In Section 4.3 is a proof for Theorem 3.2. Afterwards in Section 4.5, a lemma is presented that is very useful to prove the correctness of the enumeration algorithm and some facts for queries with aggregates. The proof of Theorem 3.3 will be continued with proving part (a) in Section 4.4, part (b) in Section 4.6, and part (c) in Section 4.7.

### 4.1. The fundamental data structure

We now give an alternative characterisation of q-hierarchical queries that sheds more light on their “tree-like” structure and will be useful for designing efficient query evaluation algorithms. Next, we define the notion of a *q-tree* for a CQ  $Q$  and show that  $Q$  is q-hierarchical if and only if it has a q-tree.

**Definition 4.1.** *Let  $Q$  be a CQ. A q-tree for  $Q$  is a rooted directed tree  $T_Q = (V, E)$  with  $V = \text{vars}(Q) \dot{\cup} \{v_{\text{root}}\}$  where*

- (1) *for all atoms  $\psi$  in  $Q$  the set  $\text{vars}(\psi) \cup \{v_{\text{root}}\}$  forms a directed path in  $T_Q$  that starts from the root  $v_{\text{root}}$ , and*
- (2) *if  $\text{free}(Q) \neq \emptyset$ , then  $\text{free}(Q) \cup \{v_{\text{root}}\}$  is a connected subset in  $T_Q$ .*

See Figure 4.1 for examples of q-trees. The following lemma gives a characterisation of the *q-hierarchical* conjunctive queries via q-trees.

**Lemma 4.2.** *A CQ  $Q$  is q-hierarchical if and only if  $Q$  has a q-tree. Moreover, there is a polynomial time algorithm that decides whether an input CQ  $Q$  is q-hierarchical, and if so, outputs a q-tree.*

*Proof.* The proof of the “only if” direction is easy, as every conjunctive query that has a q-tree  $T$  must be q-hierarchical, because if  $y$  is a descendant of  $x$  in  $T$ , then  $\text{atoms}(y) \subseteq \text{atoms}(x)$  and if the path from  $v_{\text{root}}$  to  $x$  and the path from  $v_{\text{root}}$  to  $y$  are disjoint, then  $\text{atoms}(y) \cap \text{atoms}(x) = \emptyset$ .

For proving the “if” direction, we inductively construct a q-tree  $T_Q$  for all conjunctive queries  $Q$  with at most  $s$  variables. The induction basis for empty queries is

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

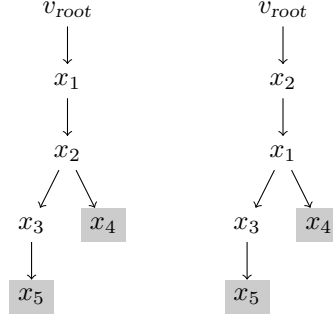


Figure 4.1.: Two  $q$ -trees for

$$Q = \{(x_1, x_2, x_3) : \exists x_4 \exists x_5 (Ex_1x_2 \wedge Rx_4x_1x_2x_1 \wedge Rx_5x_3x_2x_1)\}.$$

trivial. We just use the tree  $T = (V, E)$  with  $V = \{v_{root}\}$  and  $E = \emptyset$ . For the inductive step, we split the set of atoms in  $Q$  that have at least one variable, into disjoint sets  $\Phi_1 \cup \dots \cup \Phi_M$  we have

- for all  $m \in [M]$  the following holds. For all  $x, y \in \text{vars}(\Phi_m)$  we have  $\text{atoms}(x) \subseteq \text{atoms}(y)$  or  $\text{atoms}(x) \supseteq \text{atoms}(y)$  and
- for all  $m, m' \in [M]$  with  $m \neq m'$  the following holds. For all  $x \in \text{vars}(\Phi_m)$  and  $y \in \text{vars}(\Phi_{m'})$  we have that  $\text{atoms}(x) \cap \text{atoms}(y) = \emptyset$ .

Here, we define  $\text{vars}(\Phi) := \bigcup_{\psi \in \Phi} \text{vars}(\psi)$  for every  $\Phi \subseteq \text{atoms}(Q)$ . It follows that for all  $m, m' \in [M]$  with  $m \neq m'$ , we have that  $\Phi_m \cap \Phi_{m'} = \emptyset$ . To prove that fact let us assume for a contradiction that there is  $m, m' \in [M]$  such that there is an atom  $\psi \in \Phi_m \cap \Phi_{m'}$ . Let  $x \in \text{vars}(\psi)$  be arbitrary. Then, there is a  $x \in \text{vars}(\Phi_m) \cap \text{vars}(\Phi_{m'})$  such that  $\psi \in \text{atoms}(x)$ . This contradicts to the second condition of the definition of  $\Phi_m$  and  $\Phi_{m'}$ .

First of all, we show that  $\Phi_m = \bigcup_{w \in \text{vars}(\Phi_m)} \text{atoms}(w)$ . The fact  $\Phi_m \subseteq \bigcup_{w \in \text{vars}(\Phi_m)} \text{atoms}(w)$  is clear. To show that  $\Phi_m \supseteq \bigcup_{w \in \text{vars}(\Phi_m)} \text{atoms}(w)$ , let us assume for a contradiction that there is a  $\psi \in \bigcup_{w \in \text{vars}(\Phi_m)} \text{atoms}(w)$  with  $\psi \notin \Phi_m$ . Then,  $\psi \in \Phi_{m'}$  for a  $m' \neq m$ . Since  $\psi \in \text{atoms}(w)$  it follows that  $w \in \text{vars}(\Phi_{m'})$ . Then, there is a  $w \in \text{vars}(\Phi_m) \cap \text{vars}(\Phi_{m'})$  with  $\psi \in \text{atoms}(w)$ . This violates the second condition in the definition of the sets  $\Phi_1, \dots, \Phi_M$ .

Now, we show that for every  $m \in [M]$  there is a  $v \in \text{vars}(\Phi_m)$  such that  $v$  occurs in every atom in  $\Phi_m$ . Let  $v \in \text{vars}(\Phi_m)$  be a variable such that  $|\text{atoms}(v)| = \max_{w \in \text{vars}(\Phi_m)} |\text{atoms}(w)|$ . Thus, for all  $w \in \text{vars}(\Phi_m) \setminus \{v\}$  it follows that  $|\text{atoms}(w)| \leq |\text{atoms}(v)|$  and therefore  $\text{atoms}(w) \not\supseteq \text{atoms}(v)$  and thus  $\text{atoms}(w) \subseteq \text{atoms}(v)$ . Because of this fact it follows that  $\Phi_m = \bigcup_{w \in \text{vars}(\Phi_m)} \text{atoms}(w) = \text{atoms}(v)$  and therefore,  $v$  appears in every atom. We show now that if there are free variables in  $\Phi_m$ , the variable  $v$  that appears in every atom in  $\Phi_m$  is a free variable in  $Q$ , i.e., if  $\text{free}(Q) \cap \text{vars}(\Phi_m) \neq \emptyset$ , there must be a  $v \in \text{free}(Q) \cap \text{vars}(\Phi_m)$  with  $\Phi_m = \text{atoms}(v)$ . For a contraction let us assume that for all  $v \in \text{vars}(\Phi_m)$  with  $\Phi_m = \text{atoms}(v)$  it holds

#### 4.1. The fundamental data structure

that  $v \notin \text{free}(Q)$ . Then, by the second condition of Definition 3.1 it follows that for all  $w \in \text{vars}(\Phi_m)$  with  $\Phi_m \neq \text{atoms}(w)$  that  $w \notin \text{free}(Q)$  (since  $\text{atoms}(w) \subset \text{atoms}(v)$ ) and therefore  $\text{free}(Q) = \emptyset$ , which violates the assumption.

If  $M = 1$ , let  $v$  be the variable that appears in every atom. Let  $Q'$  be the query we obtain if we delete every atom of the form  $Rv$  and delete every appearance of  $v$  in every atom. Note that by definition  $Q'$  is q-hierarchical with  $s - 1$  variables and let  $T'$  be the querytree of  $Q'$ , we obtain from the induction hypothesis. Let  $T$  be the tree we obtain from  $T'$  if we replace the  $v_{\text{root}}$  node in  $T'$  by  $v$  and insert a “new root”  $v_{\text{root}}$  and an edge from  $v_{\text{root}}$  to  $v$ . Then,  $T$  is a q-tree for  $Q$  since for every atom  $\psi \in \text{atoms}(Q')$  the set  $(\text{vars}(\psi) \setminus \{v\}) \cup \{v_{\text{root}}\}$  forms a path in  $T'$  and by construction of  $T$  we have that the set  $(\text{vars}(\psi) \setminus \{v\}) \cup \{v\} \cup \{v_{\text{root}}\}$  forms a path in  $T$  for every  $\psi \in \text{atoms}(Q)$ . If  $\text{free}(Q) \neq \emptyset$ , then  $v \in \text{free}(Q)$  and therefore by the induction hypothesis it follows that  $\text{free}(Q) \cup \{v_{\text{root}}\}$  is a connected subset in  $T$ .

If  $M > 1$ , for every  $m \in [M]$  let  $Q_m$  be the query we obtain, by removing every atom in  $Q$  that does not belong to  $\Phi_m$  and the variables that do not belong to  $\text{vars}(\Phi_m)$ . Note that  $|\text{vars}(Q_m)| = |\text{vars}(\Phi_m)| < s$ . Thus, let  $T_m$  be the q-tree we obtain from the induction hypothesis. For every  $m \in [M]$  there must be a node  $v_m$  in  $T_m$  such that this is the only node that has an edge from the root node  $v_{\text{root}}$ . Since otherwise, if there are two nodes  $v_m^1$  and  $v_m^2$  with  $v_m^1 \neq v_m^2$  and  $\{v_{\text{root}}, v_m^1\}, \{v_{\text{root}}, v_m^2\} \in E(T_m)$ , there is an atom  $\psi_1 \in \text{atoms}(v_m^1)$  and  $\psi_2 \in \text{atoms}(v_m^2)$ , but  $\text{vars}(\psi_1) \cap \text{vars}(\psi_2) = \emptyset$ , which is a contradiction to the fact that there must be a variable that appears in every atom. Let  $T$  be the tree that we obtain from the q-trees  $T_1, \dots, T_M$  if we remove the  $v_{\text{root}}$  nodes in every tree and insert a new root node  $v_{\text{root}}$  and for all  $m \in [M]$  the edge  $\{v_{\text{root}}, v_m\}$ , i.e.,  $T = (V, E)$  is the q-tree with  $V := \bigcup_{m=1}^M V_m$  and  $E := \bigcup_{m=1}^M E_m$  where  $(V_m, E_m) = T_m$  for all  $m \in [M]$ .  $T$  is a q-tree for  $Q$  since for every atom  $\psi \in \text{atoms}(Q)$  there is a  $m \in [M]$  such that  $\psi \in \Phi_m$  and  $\text{vars}(\psi) \cup \{v_{\text{root}}\}$  forms a path in  $T_m$ . By construction, the path still appears in  $T$ . Furthermore, by the induction hypothesis it follows that  $\text{free}(Q) \cup \{v_{\text{root}}\}$  is a connected subset of in  $T$ . Thus,  $T$  is a q-tree for  $Q$ .

It is easy to see that this construction can be computed in polynomial time.  $\square$

For the remainder of this chapter we assume that  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$  is a q-hierarchical conjunctive query,  $\text{vars}(Q) = \{x_1, \dots, x_m\}$  with  $0 \leq k \leq m$ , and  $Q$  is of the form

$$Q = \{(x_1 \dots x_k, b_1, \dots, b_\ell) : \exists x_{k+1} \dots \exists x_m (\psi_1 \wedge \dots \wedge \psi_d)\}, \quad (4.1)$$

where  $b_1, \dots, b_\ell \in \mathbf{dom}$  and  $\psi_1, \dots, \psi_d$  are atomic queries of schema  $\sigma$ . We know that  $Q$  has a q-tree (see Lemma 4.2). We use the lemma's algorithm to construct in time  $\text{poly}(Q)$  a q-tree  $T_Q$  of  $Q$ . For the remainder of this section, we simply write  $T$  to denote  $T_Q$ . Recall that the vertex set  $V$  of  $T$  is the set of variables in  $Q$  and  $v_{\text{root}}$ , i.e.,  $V = \{x_1, \dots, x_m, v_{\text{root}}\}$ .

The following notation will be convenient for describing and analysing our algorithms. For a node  $v$  of  $T$ , we write  $\text{path}[v]$  to denote the set of all nodes of  $T$  that occur in the path from  $v_{\text{root}}$  to  $v$  in  $T$  (including  $v$ ), and we let  $\text{path}[v] := \text{path}[v] \setminus \{v\}$ . Furthermore we let  $\text{vpath}[v] := \text{path}[v] \setminus \{v_{\text{root}}\}$  and  $\text{vpath}[v] := \text{path}[v] \setminus \{v_{\text{root}}\}$ .

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

$child(v) := \{u : (v, u) \in E(T)\}$  is the set of children of  $v$  in  $T$ . The set  $succ(v)$ , the successors of  $v$ , is the set of variables  $w$  in  $T$  (including  $v$ ) such that there is a path from  $v$  to  $w$  in  $T$ .

A node  $v$  of  $T$  *represents* an atomic query  $\psi$  if and only if  $vpath[v] = vars(\psi)$ , i.e., the variables in  $vpath[v]$  are exactly the variables in  $\psi$ .

An *assignment* is a partial mapping from **var** to **dom**. As usual, we write  $dom(\alpha)$  for the domain of  $\alpha$ . For a set  $S \subseteq \mathbf{var}$ , by  $\alpha|_S$  we denote the restriction of  $\alpha$  to  $dom(\alpha) \cap S$ . For  $x \in \mathbf{var}$  and  $a \in \mathbf{dom}$  we write  $\alpha_x^a$  for the assignment  $\alpha'$  with domain  $dom(\alpha) \cup \{x\}$ , where  $\alpha'(x) = a$  and  $\alpha'(y) = \alpha(y)$  for all  $y \in dom(\alpha) \setminus \{x\}$ . An assignment  $\beta$  is called an *expansion* of  $\alpha$  (for short:  $\beta \supseteq \alpha$ ) if  $dom(\beta) \supseteq dom(\alpha)$  and  $\beta|_{dom(\alpha)} = \alpha$ . The *empty assignment*  $\emptyset$  is the assignment with empty domain. For pairwise distinct variables  $v_1, \dots, v_s$  and constants  $a_1, \dots, a_s \in \mathbf{dom}$  we write  $\frac{a_1, \dots, a_s}{v_1, \dots, v_s}$  to denote the assignment  $\alpha$  with  $dom(\alpha) = \{v_1, \dots, v_s\}$  and  $\alpha(v_i) = a_i$  for all  $i \in [s]$ . With  $\gamma_{id}$  we denote the function with  $dom(\gamma_{id}) = \mathbf{dom}$  and  $\gamma_{id}(a) = a$  for all  $a \in \mathbf{dom}$ . For every conjunctive formula  $\varphi$  and for every assignment  $\alpha$  with  $dom(\alpha) = vars(\varphi)$  we write  $(D, \alpha) \models \varphi$  if and only if for the valuation  $\beta = \alpha \cup \gamma_{id}$  holds  $(D, \beta) \models \{(y_1, \dots, y_s) : \varphi\}$  where  $\{y_1, \dots, y_s\} = free(\varphi)$ .

We now describe the data structure  $D$  that will be built by the **preprocess** routine and maintained while executing the **update** routine. Our data structure for a given database  $D$  represents so-called *items*. Each item  $i$  is determined by an assignment  $\alpha^i : vpath[v^i] \rightarrow \mathbf{dom}$  for some  $v^i \in V$ . We call  $\alpha^i$  the *assignment of item  $i$* , the last vertex  $v^i$  on the path the *variable of item  $i$*  and its value  $a^i := \alpha^i(v^i)$  the *constant of item  $i$* . An item  $i$  is denoted by  $[\alpha^i] = [\frac{a_1, \dots, a_s}{v_1, \dots, v_s}]$ , where  $v_{root}, v_1, \dots, v_s = v^i$  is the path from  $v_{root}$  to  $v^i$  in  $T$ . For every item  $i$  and every child  $u$  of  $v^i$  in  $T$  there is a doubly linked list  $\mathcal{L}_u^i$  (the  $u$ -list of  $i$ ) that contains items of the form  $[\alpha^i \frac{b}{u}]$  and we have one pointer  $child_u^i$  that points from  $i$  to the first element in  $\mathcal{L}_u^i$ . Note that in every data structure there is a designated *start-item*  $i_{start}$ , denoted by  $[\emptyset]$ , with the assignment  $\alpha^{i_{start}} = \emptyset$  and  $v^{i_{start}} = v_{root}$ .

It is important to note that not every item of the form  $[\alpha^i \frac{b}{u}]$  that is present in our data structure will be contained in the corresponding list  $\mathcal{L}_u^i$ . The *parent item* of an item  $i = [\frac{a_1, \dots, a_s}{v_1, \dots, v_s}]$  is defined to be the item  $[\frac{a_1, \dots, a_{s-1}}{v_1, \dots, v_{s-1}}]$ . The *start* item has no parent item. As a convention, we set  $atoms(v_{root}) = atoms(Q)$ .

Let us now state which items are actually contained in our data structure. An item  $i$  is present in our data structure if and only if there is an atom  $\psi \in atoms(v^i)$  such that there is an expansion  $\beta \supseteq \alpha^i$  with  $dom(\beta) = vars(\psi)$  and  $(D, \beta) \models \psi$ . It follows that every fact  $R(a_1, \dots, a_r)$  in the database gives rise to a constant number of items and that the overall number of items in our data structure is therefore linear in the size of the database. The definition also ensures that whenever an item is present in our data structure, then so is its parent item.

Let us now specify which of the present items are actually contained in the corresponding list  $\mathcal{L}_u^i$ .

**Definition 4.3** (fit items). *An item  $i$  is fit if and only if there is an expansion  $\beta \supseteq \alpha^i$  such that  $(D, \beta) \models \bigwedge_{\psi \in atoms(v^i)} \psi$ .*



#### 4.1. The fundamental data structure

The doubly linked list  $\mathcal{L}_u^i$  contains precisely those items that are fit. Being fit is a necessary requirement for participating in the query result and this is why we exclude unfit items from the lists. Note that whenever a tuple is inserted into or deleted from the database, this affects the “fit”-status of only a constant number of items. Furthermore, provided we have constant-time access to the items, we can update their status and include or exclude them from the corresponding lists in constant time.

The height of a variable  $v$  in a q-tree  $T_Q$  (or in a subtree  $T'$  of  $T_Q$ ) is defined as the number of edges of the longest path from  $v$  to a leaf in  $T_Q$  (or in  $T'$ , resp.). The height of an item  $i$  is defined as the height of its variable in  $T_Q$  (or in  $T'$ , resp.). This will be convenient to prove facts via induction.

We say that a data structure is lexicographically ordered if for all present items  $i$  we have that for all  $u \in \text{child}(v^i)$  the following holds. The items in the  $u$ -list of  $i$  are lexicographically ordered. This is a necessary requirement to enumerate the tuples in the query result in lexicographical order.

We consider now an example:

**Example 4.4.** *Let us consider the query*

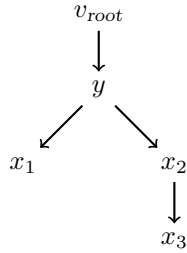
$$Q = \{(y, x_1, x_2, x_3) : E y x_1 \wedge F y x_2 x_3 \wedge G y x_2 x_3\}$$

and the database  $D_0$  with **dom** =  $\mathbb{N}$  and

$$\begin{aligned} E^{D_0} &:= \{(1, 1), (1, 2), (1, 3), (2, 4), (2, 8), (2, 9), (3, 2)\}, \\ F^{D_0} &:= \{(1, 4, 1), (1, 5, 2), (1, 6, 3), (1, 6, 4), (2, 2, 1), (2, 2, 8), (2, 2, 4), \\ &\quad (3, 1, 1), (4, 5, 6)\}, \\ G^{D_0} &:= F^{D_0}. \end{aligned}$$

See Figure 4.2 for a q-tree of  $Q$  and Figure 4.3 for an illustration of the data structure of  $Q$  on  $D$ . An arrow from an item  $i$  to an item  $i'$  labelled by  $u$  denotes that  $i'$  is the first item of the  $u$ -list of  $i$ . A line between two items  $i$  (on the left) and  $i'$  (on the right) denotes that  $i'$  is the successor item of  $i$  in the doubly-linked list. The item  $[y]$  is present in the data structure but not fit since  $(4, 5, 6) \in F^{D_0}$  and there is no  $a \in \text{adom}(D)$  with  $(4, a) \in E^{D_0}$ . All the other present items are fit.

Figure 4.2.: q-tree of  $T_Q$  from Example 4.4.





#### 4.2. Preprocessing and updating the data structure

Figure 4.4.: Illustration of the data structure for the query and database in Example 4.4 after  $(4, 1)$  was inserted to  $E^D$ .

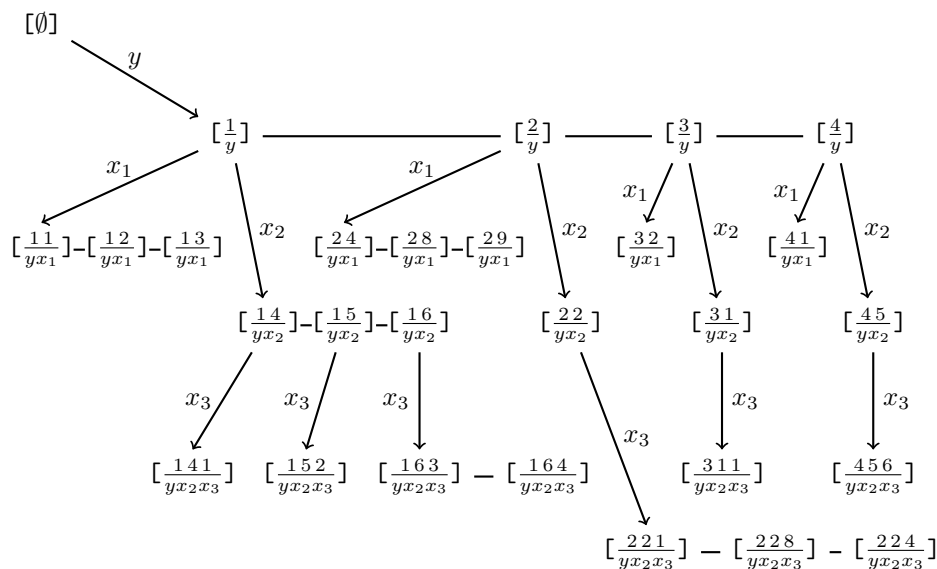
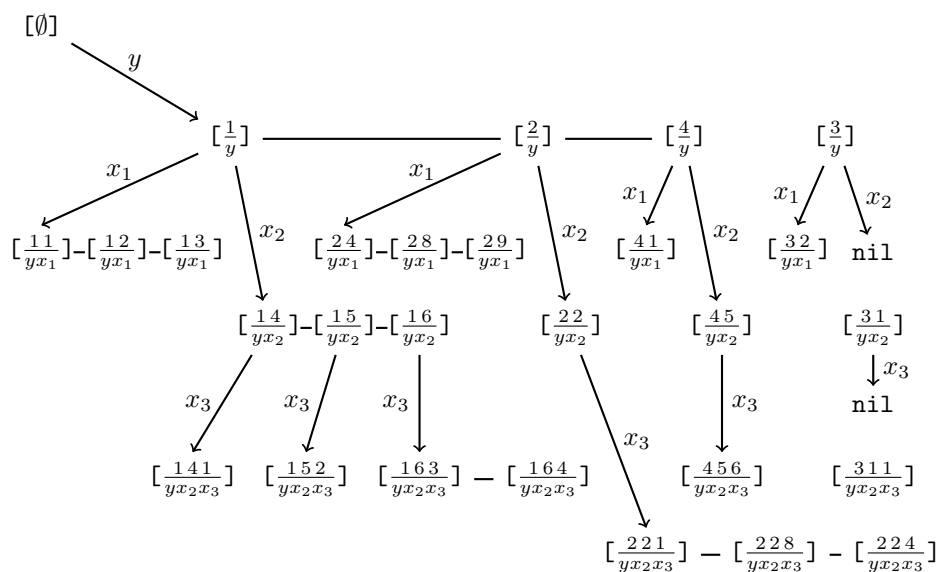


Figure 4.5.: Illustration of the data structure for the query and database in Example 4.4 after  $(4, 1)$  was inserted to  $E^D$  and  $(3, 1, 1)$  deleted from  $G^D$ .



#### 4. Answering $q$ -hierarchical conjunctive queries under updates

When an update command arrives, we have to modify our data structure accordingly so that it meets the requirements described above. For convenience, we summarise the conditions below.

**Condition 4.5.** (a) *An item  $i$  is present in our data structure if and only if there is an atom  $\psi \in \text{atoms}(v^i)$  such that there is an expansion  $\beta \supseteq \alpha^i$  with  $\text{dom}(\beta) = \text{vars}(\psi)$  and  $(D, \beta) \models \psi$ .*

(b) *For every item  $i$  and every  $u \in \text{child}(v^i)$  the list  $\mathcal{L}_u^i$  consists of all fit items of the form  $[\alpha^i \frac{b}{u}]$ .*

The update procedure that guarantees that Condition 4.5 holds is given in Algorithm 1.

We also store for every  $v \in \mathbf{var}(Q)$  a dictionary  $\mathcal{D}_v$ . Let  $v_{\text{root}}, u_1, \dots, u_s = v$  be the path from the root to  $v$  in the  $q$ -tree of  $Q$ . As keys we use  $s$ -1-ary tuples and the values are lists. In these dictionaries we represent information that an item is present, i.e.,

$$a_s \text{ is in } \mathcal{D}_{u_s}[a_1, \dots, a_{s-1}] \iff [\frac{a_1, \dots, a_s}{u_1, \dots, u_s}] \text{ is present in the data structure}$$

We additionally store for every  $R \in \sigma$  a dictionary  $\mathcal{D}_R$  with  $\text{ar}(R)$ -ary keys and  $\mathcal{D}_R[a_1, \dots, a_r] = 1$  if and only if  $(a_1, \dots, a_r) \in R^D$  ( $\mathcal{D}_R[a_1, \dots, a_r] = 0$  otherwise). Furthermore we have to store a Boolean variable **start-is-fit** that is set to **True** if, and only if, the *start*-item  $[\emptyset]$  is fit. We will later use (for example in the enumeration algorithm) this variable to check in constant time if the *start*-item is fit. Algorithm 1 shows how to change the data structure after the databases is updated, i.e., a tuple is inserted or deleted.

The correctness of Algorithm 1 follows from the following claim:

**Claim 4.6.** (a) *Condition 4.5 holds for the data structure for the empty database.*

(b) *If Condition 4.5 holds for the data structure, the condition still holds after an update was performed.*

*Proof.* (a) For the empty database we do not have items except the *start*-item and therefore the condition holds.

(b) Let us assume that we have a database  $D$  and a data structure for  $D$  and  $Q$  for that Condition 4.5 holds and we receive an update and use Algorithm 1 to update the data structure. Let  $D'$  be the resulting database after the update was performed, i.e.,  $R^{D'} = R^D \cup \{(a_1, \dots, a_r)\}$  and  $S^{D'} = S^D$  for all  $S \in \sigma \setminus \{R\}$  if we insert  $(a_1, \dots, a_r)$  to  $R$  and  $R^{D'} = R^D \setminus \{(a_1, \dots, a_r)\}$  and  $S^{D'} = S^D$  for all  $S \in \sigma \setminus \{R\}$  if we delete  $(a_1, \dots, a_r)$  from  $R$ .

Let us first consider the case that the if conditions in 8 and 10 do not hold.

The only items  $i$  that we consider in the algorithm are the following: There is an atom of the form  $Rv_1, \dots, v_r$  in  $Q$  and a  $j \in [s]$  such that  $\alpha^i = \frac{c_1, \dots, c_j}{u_1, \dots, u_j}$ . Here

## 4.2. Preprocessing and updating the data structure

---

**Algorithm 1** Procedure that modifies the data structure after an insert or delete command receives.

---

```

1: procedure UPDATE R( $a_1, \dots, a_r$ )
2:   Set  $\mathcal{D}_R[a_1, \dots, a_r] = 1$  if UPDATE == insert.
3:   Set  $\mathcal{D}_R[a_1, \dots, a_r] = 0$  if UPDATE == delete.
4:   for all atoms of the form  $Rv_1 \dots v_r$  in  $Q$  do
5:     Let  $u_1, \dots, u_s$  be the nodes in a path in the q-tree  $T_Q$  such that
        $\{u_1, \dots, u_s\} = \{v_1, \dots, v_r\} \cap \mathbf{var}$ .
6:     Let  $c_1, \dots, c_s$  be elements such that  $c_j = a_m$  if and only if  $u_j = v_m$ .
7:      $\triangleright$  The next if-condition checks if the tuple matches to constants
8:     if there is an  $i \in [r]$  with  $v_i \in \mathbf{dom}$  and  $v_i \neq a_i$  then
9:       continue  $\triangleright$  go to the next iteration step in the for loop
10:    if there are  $i, j \in [r]$  with  $v_i = v_j$  and  $u_i \neq u_j$  then
11:      continue  $\triangleright$  go to the next iteration step in the for loop
12:    if UPDATE == insert then
13:      for  $j = 1$  to  $s$  do
14:        if  $[\frac{c_1, \dots, c_j}{u_1, \dots, u_j}]$  does not exists in the data structure then
15:          Create  $[\frac{c_1, \dots, c_j}{u_1, \dots, u_j}]$  and add  $c_j$  to  $\mathcal{D}_{u_j}[c_1, \dots, c_{j-1}]$ .
16:      for  $j = s$  to  $1$  do
17:        if  $[\frac{c_1, \dots, c_j}{u_1, \dots, u_j}]$  becomes fit, add the item to  $x_j$ -list of  $[\frac{c_1, \dots, c_{j-1}}{u_1, \dots, u_{j-1}}]$ .
18:    if UPDATE == delete then
19:      for  $j = s$  to  $1$  do
20:        if  $[\frac{c_1, \dots, c_j}{u_1, \dots, u_j}]$  is not fit then
21:          Remove the item from the  $u_j$ -list of  $[\frac{c_1, \dots, c_{j-1}}{u_1, \dots, u_{j-1}}]$ .
22:        if Condition 4.5(a) does not hold for  $[\frac{c_1, \dots, c_j}{u_1, \dots, u_j}]$  then
23:          Remove  $[\frac{c_1, \dots, c_j}{u_1, \dots, u_j}]$  and remove  $c_j$  from  $\mathcal{D}_{u_j}[c_1, \dots, c_{j-1}]$ .
24:  Set start-is-fit to True if and only if  $[\emptyset]$  is fit.

```

---

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

and in the rest of the proof,  $c_1, \dots, c_s, u_1, \dots, u_s$  are defined as in lines 5 and 6 in Algorithm 1. The corresponding atom in the iteration of the for-loop should be clear from context.

If we *insert* a tuple to a relation, we create all these items  $[\frac{a_1, \dots, a_j}{u_1, \dots, u_j}]$ , since  $\alpha^i \subseteq \{v_\ell \mapsto a_\ell : \ell \in [r]\} =: \beta$  and it holds that  $(D', \beta) \models \psi$  for an atom  $\psi \in \text{atoms}(v^i)$ , so they have to be present. We also check if they are fit and add them to the corresponding lists. If we *delete* a tuple from a relation, the algorithms checks for these items if they are still fit after the update and removes them from the corresponding lists if not and removes them from the data structure if they are not allowed to be present. Note that by definition of being fit or present such items can only receive the status of being present or fit if we insert a tuple, since we insert a tuple to  $D$  and lose their status of being present or fit after we delete a tuple from a relation.

Note that, for the other items  $i$ , that we do not consider in the algorithm, it holds that for all atoms of the form  $Rv_1, \dots, v_r$  there is at least one  $\{x \mapsto b\} \in \alpha^i$  such that  $x \notin \{v_1, \dots, v_r\}$  or  $b \notin \{a_1, \dots, a_r\}$ . Furthermore, for these items  $i$  it holds that for every atom  $\psi \in \text{atoms}(v^i)$  and for every expansion  $\beta \supseteq \alpha^i$  with  $\text{dom}(\beta) = \text{vars}(\psi)$  holds  $(D, \beta) \models \psi \iff (D', \beta) \models \psi$  (and, in particular,  $(D, \beta) \models \bigwedge_{\psi \in \text{atoms}(v^i)} \psi \iff (D', \beta) \models \bigwedge_{\psi \in \text{atoms}(v^i)} \psi$ ), since the update does not change relevant tuples for their satisfaction in the database. Thus, we obtain that these items do not change the status of being present and fit.

Let us now consider the case that at least one of the if conditions in Line 8 and 10 hold. Then, for the assignment  $\beta$  with  $\beta(u_j) = c_j$  for all  $j \in [s]$  it does not hold that  $(D, \beta) \models Rv_1 \dots v_r$  and therefore such an update in consideration of the atomic formula  $Rv_1 \dots v_r$  has no affect of the present status of any item in the data structure.

Therefore, Algorithm 1 guarantees that Condition 4.5 holds after an update was performed.  $\square$

In the remainder of this section we show that the algorithm takes time  $\text{poly}(Q)$ . For all  $v \in V(T)$  let  $\text{exatoms}(v) := \left( \text{atoms}(v) \setminus \bigcup_{u \in \text{child}(v)} \text{atoms}(u) \right)$  if  $v \neq v_{\text{root}}$  and  $\text{exatoms}(v_{\text{root}})$  is the set of atoms with arity 0 in  $Q$ . First of all, we show how to quickly check that if an item is fit or not and whether an item is allowed to be present.

To do this, we use the following lemmas:

**Lemma 4.7.** *Let  $Q$  be a  $q$ -hierarchical CQ and  $D$  be a database. For all items  $i$  in the data structure for  $Q$  on  $D$  the following holds. The item  $i$  is fit if and only if for every  $u \in \text{child}(v^i)$  is the  $u$ -list of  $i$  not empty and  $(D, \alpha^i) \models \bigwedge_{\psi \in \text{exatoms}(v^i)} \psi$ .*

The proof can be found at the end of this section (see page 39). To check if an item  $i$  is fit, we simply use the fact of Lemma 4.7 and check if there is a  $u \in \text{child}(v^i)$  such that the  $u$ -list of  $i$  is empty and  $(D, \alpha^i) \models \bigwedge_{\psi \in \text{exatoms}(v^i)} \psi$ . To verify if  $(D, \alpha^i) \models \bigwedge_{\psi \in \text{exatoms}(v^i)} \psi$  we lookup for every  $Rx_1, \dots, x_r \in \text{exatoms}(v^i)$  if  $(\beta(x_1), \dots, \beta(x_r))$

#### 4.2. Preprocessing and updating the data structure

is in  $\mathcal{D}_R$  for  $\beta := \alpha^i \cup \gamma_{\text{id}}$ . Note that checking if an item is fit takes  $\text{poly}(Q)$  time, since we have at most  $|Q|$  lookups and we check for  $|\text{child}(v^i)| \leq \text{poly}(Q)$  lists, whether they are empty. Note that we update the items in the data structure bottom up, i.e., every time we check if an  $i$  is fit, the items whose assignment that are extensions of  $\alpha^i$ , already have an updated status (the information if they are fit or present).

To check if an item has to be present, we use the following lemma.

**Lemma 4.8.** *Let  $Q$  be a  $q$ -hierarchical CQ and  $D$  be a database. For all items  $i$  in the data structure for  $Q$  on  $D$  it holds that  $i$  is present if and only if there is an atom  $\psi \in \text{exatoms}(v^i)$  such that  $(D, \alpha^i) \models \psi$  or there is a  $u \in \text{child}(v^i)$  and a  $b \in \text{adom}(D)$  such that there is a present item  $\hat{i} = [\alpha^i \frac{b}{u}]$ .*

The proof is given at the end of this section (see page 40).

To check if we have to delete an item from our data structure in line 22 of Algorithm 1, we check if

1. there is an atom  $\psi \in \text{exatoms}(v^i)$  such that  $(D, \alpha^i) \models \psi$  or
2. there is a  $u \in \text{child}(v^i)$  and a  $b \in \text{adom}(D)$  such that there is a present item  $\hat{i} = [\alpha^i \frac{b}{u}]$ .

To check if condition 1 holds we lookup for all  $\psi = Rx_1 \dots x_r \in \text{exatoms}(v^i)$  if  $(\beta(x_1), \dots, \beta(x_r))$  is in  $\mathcal{D}_R$  where  $\beta = \alpha^i \cup \gamma_{\text{id}}$ . For condition 2, we check if  $\mathcal{D}_{x_r}[\alpha(x_1), \dots, \alpha(x_{r-1})]$  is not empty.

Since lines 8, 10, 17, 20, 22 and 24 in Algorithm 1 take time  $\text{poly}(Q)$  and the for loop has at most  $\|Q\|$  iterations and the other lines take constant time, we need time  $\text{poly}(Q)$  for Algorithm 1. Thus, we can update the data structure, that is not lexicographically ordered, in time  $\text{poly}(Q)$  if we receive an update.

**Remark 4.9.** *Note that for lexicographically ordered data structures we have to take into account that for all present item  $i$  in the data structure and for all  $u \in \text{child}(v^i)$  in the data structure the  $u$ -list of  $i$  is lexicographically ordered. To realise this, we have to create for all lists ordered data structures such as AVL-trees or red-black trees (see [32] for a survey). For these data structures it takes time  $O(\log(|\text{adom}(D)|))$  to insert or to delete an element from the list. By analysing the running time for the update algorithm with the assumption that it takes time  $O(\log(|\text{adom}(D)|))$  to insert or delete an element from a  $u$ -list of  $i$ , we obtain that it takes time  $\log(|\text{adom}(D)|) \text{poly}(Q)$  to update lexicographically ordered data structures.*

In the remainder of this section we give proofs for Lemma 4.7 and Lemma 4.8.

*Proof of Lemma 4.7.* Let us assume that a present item  $i$  is fit. Then, there is a  $\beta \supseteq \alpha^i$  such that  $(D, \beta) \models \bigwedge_{\psi \in \text{atoms}(v^i)} \psi$ . In particular, it holds that  $(D, \alpha^i) \models \bigwedge_{\psi \in \text{exatoms}(v^i)} \psi$  and for every  $w \in \text{child}(v^i)$  it holds that  $\text{atoms}(w) \subseteq \text{atoms}(v^i)$  and therefore  $(D, \beta) \models \bigwedge_{\psi \in \text{atoms}(w)} \psi$ . Thus, there is an assignment  $\alpha_w = \beta|_{\text{vpath}[w]} \subseteq \beta$  such that  $(D, \beta) \models \bigwedge_{\psi \in \text{atoms}(w)} \psi$  and, in particular, the item  $[\alpha_w]$  is present and fit. Therefore, there must be an item  $[\alpha_w]$  in the  $w$ -list of  $i$  for every  $w \in \text{child}(v^i)$ .

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

For the “only if” direction let us assume that every  $u$ -list of  $i$  is not empty and  $(D, \alpha^i) \models \psi$  for all  $\psi \in \text{exatoms}(v^i)$ . Then, there exists for every  $u \in \text{child}(v^i)$  an item  $[\alpha^i \frac{b}{u}]$  that is fit, i.e., there is a  $\beta_u \supseteq \alpha^i \frac{b}{u}$  such that  $(D, \beta_u) \models \bigwedge_{\psi \in \text{atoms}(u)} \psi$ . Let  $\beta := \bigcup_{u \in \text{child}(v^i)} \beta_u$ . Note that  $\beta$  is a valid assignment since  $\alpha^i \subseteq \beta_u$  for all  $u \in \text{child}(v^i)$  and  $\text{dom}(\alpha^i) = \bigcap_{u \in \text{child}(v^i)} \text{dom}(\beta_u)$  and  $\text{atoms}(v^i) = \bigcup_{u \in \text{child}(v^i)} \text{atoms}(u) \cup \text{exatoms}(v)$ . It follows that  $(D, \beta) \models \bigwedge_{\psi \in \text{atoms}(v^i)} \psi$  and  $i$  is therefore fit.  $\square$

*Proof of Lemma 4.8.* Let us assume that  $i$  is a present item in the data structure and let  $\psi \in \text{atoms}(v^i)$  be the atom and  $\beta \supseteq \alpha^i$  be the assignment with  $\text{dom}(\beta) = \text{vars}(\psi)$  and  $(D, \beta) \models \psi$ .

- *Case 1:*  $\psi \in \text{exatoms}(v^i)$ . Then,  $\text{dom}(\beta) = \text{vars}(\psi) = \text{dom}(\alpha^i)$  and thus,  $(D, \alpha^i) \models \psi$ .
- *Case 2:*  $\psi \notin \text{exatoms}(v^i)$ . Then, there exists a  $u \in \text{child}(v^i)$  such that  $\psi \in \text{atoms}(u)$ . Let  $b = \beta(u)$ . Then, it holds that  $\beta \supseteq \alpha^i \frac{b}{u}$  with  $\text{dom}(\beta) = \text{vars}(\psi)$  and  $(D, \beta) \models \psi$  and, in particular, there is a present item  $[\alpha^i \frac{b}{u}]$  in the data structure.

For the “only if” direction, we consider two cases:

- *Case 1:* There is an atom  $\psi \in \text{exatoms}(v^i)$  such that  $(D, \alpha^i) \models \psi$ . It is straightforward to verify that  $i$  is present in the data structure.
- *Case 2:*  $[\alpha^i \frac{b}{u}]$  is present in our data structure. Then, there is an atom  $\psi \in \text{atoms}(u) \subseteq \text{atoms}(v^i)$  and an expansion  $\beta \supseteq \alpha^i \frac{b}{u} \supseteq \alpha^i$  with  $\text{dom}(\beta) = \text{vars}(\psi)$  and  $(D, \beta) \models \psi$ . By definition,  $i$  is present in the data structure.  $\square$

### 4.3. Answering Boolean $q$ -hierarchical queries

The aim of this section is to give a proof for Theorem 3.2, i.e., to show how to use the data structure for answering Boolean  $q$ -hierarchical conjunctive queries.

For convenience, Theorem 3.2 is restated in the following.

**Theorem 3.2.** *There is a dynamic algorithm that receives a Boolean  $q$ -hierarchical conjunctive query  $Q$  and a  $\sigma$ -db  $D_0$ , and computes within  $t_i = \text{poly}(Q)$  initialisation time and  $t_p = \text{poly}(Q) \cdot O(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)$  and allows to answer  $Q(D)$  in time  $t_{\text{ans}} = O(1)$  where  $D$  is the current database.*

*Proof of Theorem 3.2.* Let

$$Q := \{(b_1, \dots, b_\ell) : \exists y_1 \dots \exists y_k \varphi\}$$

be the input Boolean  $q$ -hierarchical CQ and let  $D$  be the input database. The main idea is to use the data structure from this chapter to maintain the query  $Q'$  on  $D$  under updates where

$$Q' := \{(y_1, \dots, y_k, b_1, \dots, b_\ell) : \varphi\}$$



with the **init**, **preprocess** and the **update** routines from Section 4.2. An update can be done in time  $t_u = \text{poly}(Q)$ , the initialisation can be done in time  $t_i = \text{poly}(Q)$  and the preprocessing can be done in time  $t_p = \|D_0\| \text{poly}(Q)$ .

For the **answer** routine, we just check if the variable **start-is-fit** is **true** and return the corresponding result. This can be done in time  $O(1)$ . The routine is correct, since by Definition 4.3 it follows that the start item is fit, if and only, if there is an assignment  $\beta$  such that  $(D, \beta) \models \bigwedge_{\psi \in \text{atoms}(v_{\text{root}})} \psi$ . Thus,  $(D, \beta) \models Q'$ . Obviously, this fact holds if and only if  $Q(D) = \text{yes}$ .

This concludes the proof of Theorem 3.2.  $\square$

## 4.4. Testing q-hierarchical queries

The aim of this section is to show Theorem 3.3 (a), i.e., to show how the test routine for q-hierarchical queries works.

First we consider an example.

**Example 4.10.** *Let us recall the database  $D_0$  and the query  $Q$  from Example 4.4. Assume, that we want to test if  $(1, 3, 6, 3)$  belongs to the result set. Let  $\alpha := \frac{1 \ 3 \ 6 \ 3}{y x_1 x_2 x_3}$  be the assignment such that it belongs to the tuple if matching the free variables with the components of the tuple, i.e.,  $(1, 3, 6, 3) = (\alpha(y), \alpha(x_1), \alpha(x_2), \alpha(x_3))$ .*

*For all  $u \in \text{free}(Q)$  we consider the assignments  $\alpha_u := \alpha|_{\text{vpath}[u]}$  and check if the items  $[\emptyset]$  and  $[\alpha_u]$  are fit for all  $u \in \text{free}(Q)$ . The following holds.*

- *Since  $[\emptyset]$  is fit, it follows that there are  $a_1, a_2, a_3, a_4 \in \text{adom}(D_0)$  such that  $(a_1, a_2) \in E^{D_0}$  and  $(a_1, a_3, a_4) \in F^{D_0} \cap G^{D_0}$ .*
- *Since  $[\frac{1}{y}]$  is fit, it follows that there are  $a_2, a_3, a_4 \in \text{adom}(D_0)$  such that  $(1, a_2) \in E^{D_0}$  and  $(1, a_3, a_4) \in F^{D_0} \cap G^{D_0}$ .*
- *Since  $[\frac{1 \ 3}{y x_1}]$  is fit, it follows that  $(1, 3) \in E^{D_0}$ .*
- *Since  $[\frac{1 \ 6}{y x_2}]$  is fit, it follows that there is a  $a_4 \in \text{adom}(D_0)$  such that  $(1, 6, a_4) \in F^{D_0} \cap G^{D_0}$ .*
- *Since  $[\frac{1 \ 6 \ 3}{y x_2 x_3}]$  is fit, it follows that  $(1, 6, 3) \in F^{D_0} \cap G^{D_0}$ .*

*Since  $(1, 3) \in E^{D_0}$  and  $(1, 6, 3) \in F^{D_0} \cap G^{D_0}$ , it follows that  $(1, 3, 6, 3)$  belongs to the query result.*

*Assume, that we want to test if  $(1, 4, 6, 3)$  belongs to the query result. There is no fit and present item  $[\frac{1 \ 4}{y x_1}]$  and thus  $(1, 4) \notin E^{D_0}$  and therefore,  $(1, 4, 6, 3)$  does not belong to the query result.*

For the remainder of this section we assume that  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$  is a q-hierarchical conjunctive query,  $\text{vars}(Q) = \{x_1, \dots, x_m\}$  with  $0 \leq k \leq m$ , and  $Q$  is of the form

$$Q = \{(x_1 \dots x_k, b_1, \dots, b_\ell) : \exists x_{k+1} \dots \exists x_m (\psi_1 \wedge \dots \wedge \psi_d)\}, \quad (4.2)$$

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

where  $b_1, \dots, b_\ell \in \mathbf{dom}$  and  $\psi_1, \dots, \psi_d$  are atomic queries of schema  $\sigma$ . We are given as input  $Q$  and a  $\sigma$ -db  $D$  and let us assume we have the data structure from Section 4.1 for  $Q$  and  $D$ . When we call the **test** routine with input  $\bar{a} = (\alpha(x_1), \dots, \alpha(x_k), c_1, \dots, c_\ell)$ , we check if the following two conditions are true.

- $b_j = c_j$  for all  $j \in [\ell]$ .
- For all  $v \in \text{vars}(Q)$  the items  $[\alpha|_{\text{vpath}[v]}]$  are fit and output the result of the test.

We output **Yes** if and only if the results of both tests are true.

Note that by Definition 4.3, it follows that

$$\begin{aligned} (D, \alpha) \models \psi \text{ for all atoms } \psi \in \text{atoms}(Q) \\ \Leftrightarrow (D, \alpha) \models \psi \text{ for all atoms } \psi \in \text{atoms}(v) \text{ for all } v \in \text{vars}(Q) \\ \text{and for all atoms } \psi \text{ with } \text{artiy } 0 \\ \Leftrightarrow [\alpha|_{\text{vpath}[v]}] \text{ is fit for all } v \in \text{vars}(Q) \end{aligned}$$

and therefore the **test** routine is correct. Checking the first condition can be done in  $O(\text{poly}(Q))$ . Using Lemma 4.7, we test for an item  $i$  if for every  $u \in \text{child}(v^i)$  the  $u$ -list of  $i$  is not empty and if  $(D, \alpha^i) \models \bigwedge_{\psi \in \text{exatoms}(v^i)} \psi$ . With the data structure in Section 4.1, this can be done in  $O(\text{poly}(Q))$  time. This has to be done  $|\text{vars}(Q)|$  times and therefore the **test** routine takes time  $O(\text{poly}(Q))$ . This concludes the proof of Theorem 3.3 (a).

### 4.5. The decomposition lemma

This section presents a lemma which is a very convenient tool to obtain a lot of results and it will be used in various manners. It gives the conditions that holds for every present and fit item in the data structure for obtaining the result set from the data structure. Already in the next Section 4.6 there is an application given, how to use this lemma to obtain the result for enumeration.

Before we give the lemma, we need a definition.

**Definition 4.11** (Extensions sets for items). *For all present items  $i$  in the data structure let*

$$\mathcal{E}^i := \left\{ \beta \supseteq \alpha^i : \text{dom}(\beta) = \bigcup_{\psi \in \text{atoms}(v^i)} \text{vars}(\psi), (D, \beta) \models \bigwedge_{\psi \in \text{atoms}(v^i)} \psi \right\}.$$

Furthermore, we let

$$\tilde{\mathcal{E}}^i := \left\{ \beta|_{\text{free}(Q)} : \beta \in \mathcal{E}^i \right\}. \quad (4.3)$$

Recall that  $\text{atoms}(v_{\text{root}}) = \text{atoms}(Q)$  and  $v^{[\emptyset]} = v_{\text{root}}$ . It follows that  $Q(D) = \{(\alpha(x_1), \dots, \alpha(x_k), b_1, \dots, b_\ell) : \alpha \in \tilde{\mathcal{E}}^{[\emptyset]}\}$  if  $[\emptyset]$  is fit. The decomposition lemma for  $q$ -hierarchical queries shows us how the sets  $\tilde{\mathcal{E}}^i$  for all fit items  $i$  in the data structure can be decomposed:

#### 4.5. The decomposition lemma

**Lemma 4.12** (Decomposition Lemma for q-hierarchical queries). *Let  $Q$  be a q-hierarchical CQ and let  $D$  be a  $\sigma$ -db and let  $T_Q$  be a q-tree of  $Q$ . Let  $i$  be a present and fit item in the data structure of  $Q$  on  $D$  and let  $u_1, \dots, u_s$  be the children of  $v^i$  in  $T_Q$ . Then,*

$$\mathcal{E}^i = \left\{ \bigcup_{j=1}^s \alpha_j \cup \alpha^i : \text{for all } j \in [s] \text{ there is an } \hat{\iota} \in \mathcal{L}_{u_j}^i \text{ such that } \alpha_j \in \mathcal{E}^{\hat{\iota}} \right\}$$

and

$$\tilde{\mathcal{E}}^i = \left\{ \bigcup_{j=1}^s \alpha_j \cup \alpha^i : \begin{array}{l} \text{for all } j \in [s] \text{ with } u_j \in \text{free}(Q) \text{ there is a } \hat{\iota} \in \mathcal{L}_{u_j}^i \\ \text{such that } \alpha_j \in \tilde{\mathcal{E}}^{\hat{\iota}} \end{array} \right\}.$$

*Proof.* Let  $\alpha \in \mathcal{E}^i$ . Then by definition,  $\alpha \supseteq \alpha^i$  is an extension with  $\text{dom}(\alpha) = \bigcup_{\psi \in \text{atoms}(v^i)} \text{vars}(\psi)$  such that, for all  $\psi \in \text{atoms}(v^i)$  we have that  $(D, \alpha) \models \psi$ . Note that by construction of the q-tree the set  $\text{atoms}(v^i)$  can be decomposed into disjoint sets  $\text{atoms}(v^i) = \text{exatoms}(v^i) \dot{\cup} \text{atoms}(u_1) \dot{\cup} \dots \dot{\cup} \text{atoms}(u_s)$ , where  $u_1, \dots, u_s$  are the children of  $v^i$  in  $T_Q$ . In particular, we obtain that for all atoms  $\psi \in \text{exatoms}(v^i)$  we have  $(D, \alpha) \models \psi$  and for all  $j \in [s]$  it holds that for all  $\psi \in \text{atoms}(u_j)$  is  $(D, \alpha) \models \psi$  since  $i$  is fit. Furthermore we can restrict the domain of the assignment  $\alpha$  to variables, that appear in the atom and we obtain for all atoms  $\psi \in \text{exatoms}(v^i)$  that  $(D, \alpha|_{\text{vpath}[v^i]}) \models \psi$  and for all  $j \in [s]$  we have that for all  $\psi \in \text{atoms}(u_j)$  it holds that  $(D, \alpha|_{\text{vpath}[v^i] \cup \text{succ}(u_j)}) \models \psi$ . By definition of fit, we obtain for all  $j \in [k]$  that the item  $[\alpha|_{\text{vpath}[u_j]}]$  is fit (since for the assignment  $\alpha|_{\text{vpath}[v^i] \cup \text{succ}(u_j)}$  it holds that  $(D, \alpha|_{\text{vpath}[v^i] \cup \text{succ}(u_j)}) \models \psi$  for all  $\psi \in \text{atoms}(u_j)$ ) and thus  $\hat{\iota} = [\alpha|_{\text{vpath}[u_j]}]$  is included in the  $u_j$ -list of  $i$ . Furthermore, by definition,  $\alpha|_{\text{vpath}[v^i] \cup \text{succ}(u_j)} \in \mathcal{E}^{\hat{\iota}}$ . We have that

$$\alpha \in \left\{ \bigcup_{j=1}^s \alpha_j \cup \alpha^i : \text{for all } j \in [s] \text{ there is a } \hat{\iota} \in \mathcal{L}_{u_j}^i \text{ such that } \alpha_j \in \mathcal{E}^{\hat{\iota}} \right\}.$$

It remains to show that every assignment on the right side of the equation belongs to  $\mathcal{E}^i$ . Let  $\alpha = \alpha_1 \cup \dots \cup \alpha_s$  be an assignment such that there is a  $\hat{\iota}_j \in \mathcal{L}_{u_j}^i$  where  $\alpha_j \in \mathcal{E}^{\hat{\iota}_j}$ . By definition, for every  $\psi \in \text{atoms}(u_j)$  it holds that  $(D, \alpha_j) \models \psi$  for all  $j \in [s]$  and, in particular,  $(D, \alpha) \models \psi$ . Since  $i$  is fit, it follows that for all  $\psi \in \text{exatoms}(v^i)$  we have  $(D, \alpha^i) \models \psi$  and, in particular,  $(D, \alpha) \models \psi$ . Since  $\text{atoms}(v^i) = \text{exatoms}(v^i) \dot{\cup} \text{atoms}(u_1) \dot{\cup} \dots \dot{\cup} \text{atoms}(u_s)$ , it holds that  $\alpha \in \mathcal{E}^i$ .

For the second equation, we use for all items  $i$  with  $v^i \in \text{free}(Q)$  the following argument.

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

$$\begin{aligned}
\alpha|_{\text{free}(Q)} \in \tilde{\mathcal{E}}^i &\iff \alpha \in \mathcal{E}^i \\
&\iff \alpha \in \left\{ \bigcup_{j=1}^s \alpha_j \cup \alpha^i : \begin{array}{l} \text{for all } j \in [s] \text{ there is a} \\ \hat{\iota} \in \mathcal{L}_{u_j}^i \text{ such that } \alpha_j \in \mathcal{E}^{\hat{\iota}} \end{array} \right\} \\
&\stackrel{(\star)}{\iff} \alpha|_{\text{free}(Q)} \in \left\{ \bigcup_{j=1}^s \alpha_j|_{\text{free}(Q)} \cup \alpha^i : \begin{array}{l} \text{for all } j \in [s] \text{ where } u_j \in \text{free}(Q) \\ \text{there is a } \hat{\iota} \in \mathcal{L}_{u_j}^i \\ \text{such that } \alpha_j|_{\text{free}(Q)} \in \tilde{\mathcal{E}}^{\hat{\iota}} \\ \text{and } \alpha_j := \alpha^i \text{ for all } j \in [s] \\ \text{with } u_j \notin \text{free}(Q) \end{array} \right\}
\end{aligned}$$

The “ $\implies$ ”-direction of  $(\star)$  is easy to see. To verify that the “ $\impliedby$ ”-direction holds let  $\beta$  be an assignment of the form  $\beta = \beta_1 \cup \dots \cup \beta_s \cup \alpha^i$  where for all  $j \in [s]$  with  $u_j \in \text{free}(Q)$  there is a  $\hat{\iota} \in \mathcal{L}_{u_j}^i$  such that  $\beta_j \in \tilde{\mathcal{E}}^{\hat{\iota}}$  and for all  $j \in [s]$  with  $u_j \notin \text{free}(Q)$  we have  $\beta_j = \alpha^i$ . By the definition of  $\tilde{\mathcal{E}}^i$  and by the fact that  $i$  is fit, we have that for all  $j \in [s]$  with  $u_j \in \text{free}(Q)$  there is a  $\hat{\iota} \in \mathcal{L}_{u_j}^i$  and a  $\alpha_j \in \mathcal{E}^{\hat{\iota}}$  where  $\alpha_j|_{\text{free}(Q)} = \beta_j$  and for all  $j \in [s]$  with  $u_j \notin \text{free}(Q)$  we have that there is a  $\hat{\iota} \in \mathcal{L}_{u_j}^i$  and a  $\beta_j \in \mathcal{E}^{\hat{\iota}}$  where  $\alpha_j|_{\text{free}(Q)} = \alpha^i$ . Let  $\alpha := \alpha_1 \cup \dots \cup \alpha_s \cup \alpha^i$ . Clearly  $\alpha|_{\text{free}(Q)} = \beta$  and

$$\alpha \in \left\{ \bigcup_{j=1}^s \alpha_j \cup \alpha^i : \begin{array}{l} \text{for all } j \in [s] \text{ there is a} \\ \hat{\iota} \in \mathcal{L}_{u_j}^i \text{ such that } \alpha_j \in \mathcal{E}^{\hat{\iota}} \end{array} \right\}.$$

□

### 4.6. Enumerating the query result with constant delay

The aim of this chapter is to prove Theorem 3.3(b), i.e. we now discuss how the data structure from the previous chapter can be used to enumerate the query result with constant delay. See Table 4.1 for the 22 result tuples of Example 4.4.

In the next example we show how to algorithmically enumerate the tuples in Table 4.1 using the data structure for  $D_0$  and  $Q$  from Example 4.4

**Example 4.13.** *Let us consider the query  $Q$  and the database  $D_0$  from Example 4.4. Algorithm 2 enumerates the tuples in the query result. The algorithm runs through the data structure and it runs over all fit items. Since all the items are fit and the output tuples comes from the assignments of fit item, it follows that the tuples are part of the query result.*

*Until we reach the first output in our example, the item in line 2 is set to  $\lceil \frac{1}{y} \rceil$ , the item in line 3 is set to  $\lceil \frac{1}{y x_1} \rceil$ , the item in line 4 is set to  $\lceil \frac{1}{y x_2} \rceil$  and the item in line 5 is set to  $\lceil \frac{1}{y x_2 x_3} \rceil$ . Since  $\lceil \frac{1}{y x_2 x_3} \rceil$  is the item in the  $x_3$ -list of  $\lceil \frac{1}{y x_2} \rceil$  and thus, in we*

#### 4.6. Enumerating the query result with constant delay

$y$	$x_1$	$x_2$	$x_3$
1	1	4	1
1	1	5	2
1	1	6	3
1	1	6	4
1	2	4	1
1	2	5	2
1	2	6	3
1	2	6	4
1	3	4	1
1	3	5	2
1	3	6	3
1	3	6	4
2	4	2	1
2	4	2	8
2	4	2	4
2	8	2	1
2	8	2	8
2	8	2	4
2	9	2	1
2	9	2	8
2	9	2	4
3	2	1	1

Table 4.1.: Enumeration of  $Q(D_0)$  from Example 4.4.

---

**Algorithm 2** Enumeration algorithm.

---

```

1: if  $[\emptyset]$  is fit then
2:   for  $[\frac{a_0}{y}]$  in the  $y$ -list of  $[\emptyset]$  do
3:     for  $[\frac{a_0 a_1}{y x_1}]$  in the  $x_1$ -list of  $[\frac{a_0}{y}]$  do
4:       for  $[\frac{a_0 a_2}{y x_2}]$  in the  $x_2$ -list of  $[\frac{a_0}{y}]$  do
5:         for  $[\frac{a_0 a_2 a_3}{y x_2 x_3}]$  in the  $x_3$ -list of  $[\frac{a_0 a_2}{y x_2}]$  do
6:           Output  $(a_0, a_1, a_2, a_3)$ 
7: Output EOE

```

---

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

set the item in line 4 is set to  $\lceil \frac{1}{y} \frac{5}{x_2} \rceil$  and in line 5 is set to  $\lceil \frac{1}{y} \frac{5}{x_2 x_3} \rceil$ . The algorithm continues with enumerating all results.

For the remainder of this chapter we assume that  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$  is a  $q$ -hierarchical conjunctive query,  $\text{vars}(Q) = \{x_1, \dots, x_m\}$  with  $0 \leq k \leq m$ , and  $Q$  is of the form

$$Q = \{(x_1 \dots x_k, b_1, \dots, b_\ell) : \exists x_{k+1} \dots \exists x_m (\psi_1 \wedge \dots \wedge \psi_d)\}, \quad (4.4)$$

where  $b_1, \dots, b_\ell \in \mathbf{dom}$  and  $\psi_1, \dots, \psi_d$  are atomic queries of schema  $\sigma$ . To enumerate the result of a non-Boolean conjunctive  $q$ -hierarchical query  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$ , let  $T_{\text{free}}$  be the subtree of  $T_Q$  induced on  $V(T_{\text{free}}) := \text{free}(Q) \cup \{v_{\text{root}}\} = \{x_1, \dots, x_k, v_{\text{root}}\}$ . For each node  $v$  of  $T_{\text{free}}$ , let us fix an (arbitrary) linear order  $\leq^v$  for all variables  $v$  in  $Q$ . In our example query we have  $T_{\text{free}} = T$  and we let  $x_1 <^y x_2$ . For the enumeration procedure we will use the notion of  $\mathcal{E}^i$  (see Definition 4.11) and the decomposition of these sets given in Lemma 4.12.

The main idea for the **enumerate** routine is the following. If the *start*-item  $[\emptyset]$  is not fit, the **enumerate** routine stops immediately with output **EOE**. Otherwise, we do the following. Inductively, for every item  $i$  we enumerate the set  $\tilde{\mathcal{E}}^i$  using the following algorithm. We iterate for all  $u \in \text{child}(v^i) \cap \text{free}(Q)$  over all  $\hat{t} \in \mathcal{L}_u^i$  over all assignments in  $\alpha_u \in \tilde{\mathcal{E}}^{i_j}$  and construct the assignment  $\alpha^i \cup \bigcup_{u \in \text{child}(v^i)} \alpha_u$ . By using Lemma 4.12 we know that we enumerate the assignments in  $\tilde{\mathcal{E}}^i$ . If we enumerate  $\tilde{\mathcal{E}}^{[\emptyset]}$ , we can easily construct the result tuples in  $Q(D)$  by outputting  $(\alpha(x_1), \dots, \alpha(x_k), b_1, \dots, b_\ell)$  for every  $\alpha \in \tilde{\mathcal{E}}^{[\emptyset]}$ . The pseudo-code for the described recursive **enumerate** routine is given in Algorithm 3. The **ENUM** function in Algorithm 3 also requests a set of an order for every node. This order defines in which grouping the result set will be output, since the order can change the sorting of the for loops. For example, in Table 4.1 we used the order  $x_1 <^y x_2$  for  $y$ . This implies by construction of the algorithm, that the result set is grouped by the variables  $y$  and  $x_1$ .

The next lemma establishes the correctness of Algorithm 3 and show that the tuples will be enumerated with delay  $O(|\text{vars}(Q)|)$  without repetition.

**Lemma 4.14.** *Let  $Q$  be a  $q$ -hierarchical CQ and let  $D$  be a  $\sigma$ -db. For all present and fit items  $i$  in the data structure for  $Q$  on  $D$  it holds that:*

- (a) *the assignments yielded by  $\text{ENUM}(i, \{<^u\}_{u \in \text{succ}(v^i)})$  in Algorithm 3 are exactly the assignments in  $\tilde{\mathcal{E}}^i$ .*
- (b) *The procedure  $\text{ENUM}(i, \{<^u\}_{u \in \text{succ}(v^i)})$  in Algorithm 3 takes time  $O(|\text{succ}(v^i)|)$* 
  - *until the first assignment will be yielded and*
  - *between two assignments were yielded and*
  - *the last assignment will be yielded and the procedure finished.*
- (c)  *$\text{ENUM}(i, \{<^u\}_{u \in \text{succ}(v^i)})$  in Algorithm 3 does not output duplicates.*

*Proof.* Let  $T_Q$  be the  $q$ -tree of  $Q$  that is used to construct the data structure for  $Q$  on  $D$  and let  $T_{\text{free}}$  be the subtree of  $T_Q$  induced on  $\text{free}(Q)$ .

#### 4.6. Enumerating the query result with constant delay

---

**Algorithm 3** Enumeration algorithm.

---

```

1: function ENUM( $i, \{<^v\}_{v \in \text{succ}(v^i)}$ )
2:   Input: present and fit item  $i$  in the data structure  $D$  and a set that contains
   for every  $w \in \text{succ}(v)$  in  $T_{\text{free}}$  an order over their children.
3:   if  $v^i$  is a leaf in  $T_{\text{free}}$  then
4:     yield  $\alpha^i$ 
5:   else
6:     Let  $u_1 <^{v^i} \dots <^{v^i} u_s$  be the children of  $v^i$  that belongs to  $\text{free}(Q)$ .
7:     for  $\hat{\ell}_1 \in \mathcal{L}_{u_1}^i$  do
8:       for  $\alpha_1 \in \text{ENUM}(\hat{\ell}_1, \{<^v\}_{v \in \text{succ}(u_1)})$  do
9:          $\vdots$ 
10:        for  $\hat{\ell}_s \in \mathcal{L}_{u_s}^i$  do
11:          for  $\alpha_s \in \text{ENUM}(\hat{\ell}_s, \{<^v\}_{v \in \text{succ}(u_s)})$  do
12:            yield  $\alpha^i \cup \alpha_1 \cup \dots \cup \alpha_s$ 
13:
14:   if  $[\emptyset]$  is fit then
15:     for  $\alpha \in \text{ENUM}([\emptyset], \{<^v\}_{v \in V'})$  do
16:       print  $(\alpha(x_1), \dots, \alpha(x_k), b_1, \dots, b_s)$ 
17: print EOE

```

---

We show Lemma 4.14 by induction of the height of an item in  $T_{\text{free}}$ .

For the induction base let us consider an item  $i$  of height 0, i.e.,  $v^i$  is a leaf in  $T_{\text{free}}$ . Then, we simply output  $\alpha^i$ . This is exactly the only assignment in  $\mathcal{E}^i$  and it takes  $O(1)$  time until the assignment will be yield and until the procedure will finish. Clearly, we do not output duplicates.

For the inductive step let us consider an item of height  $h$ . Let  $u_1, \dots, u_s \subseteq \text{free}(Q)$  be the children of  $v^i$  in  $T_{\text{free}}$ . By construction of the algorithm it holds that for all  $j \in [s]$  there is a  $\hat{\ell}_j \in \mathcal{L}_{u_j}^i$  such that  $\alpha^i \cup \alpha_1 \cup \dots \cup \alpha_s$  will be yield, where  $\alpha_j$  will be yielded by  $\text{ENUM}(\hat{\ell}_j)$ . Note that the items  $\hat{\ell}_j$  for all  $j \in [k]$  are fit. Therefore, by induction hypothesis, these are exactly the assignments in  $\tilde{\mathcal{E}}^{\hat{\ell}_j}$ . Thus, with Lemma 4.12 it follows that the assignments, yielded by  $\text{ENUM}(i, \{<^u\}_{u \in \text{succ}(v^i)})$  are exactly the assignments in  $\tilde{\mathcal{E}}^i$ .

Part (b) follows from the fact that for every iteration of the for-loop, where we iterate over assignments in  $\text{ENUM}(\hat{\ell}_j, \{<^w\}_{w \in \text{succ}(u_j)})$ , it takes by induction hypothesis time  $O(|\text{succ}(u)|)$  to receive the first assignment we consider, and between two considered assignments, and between the last assignment and the end of the for-loop. Furthermore, we need time  $O(1)$  to go to the first or the next element of  $\mathcal{L}_{u_j}^i$ . Therefore by induction hypothesis, we need  $O(\sum_{j=1}^s |\text{succ}(u_j)|) = O(\text{succ}(v^i))$  until the first assignment will be yielded and between two assignments were yielded and between the last assignment will be yielded and the end of the procedure.

To prove part (c), let us assume for a contradiction, that an assignment  $\alpha$  was yielded twice. We consider now the outermost loop, that continued with the iteration

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

between the two times  $\alpha$  was yielded. If the loop is of the form  $\widehat{t} \in \mathcal{L}_{u_j}^i$ , item  $[\alpha^{\frac{\alpha(u_j)}{u_j}}]$  was considered twice. This is a contradiction to the construction of the data structure. If the loop is of the form  $\alpha_j \in \text{Enum}(\widehat{t}_j)$ , then  $\alpha_j$  was considered twice but this is a contradiction to the induction hypothesis.  $\square$

Algorithm 3 is correct, since

$$Q(D) = \{(\alpha(x_1), \dots, \alpha(x_k), b_1, \dots, b_\ell) : \alpha \in \widetilde{\mathcal{E}}^{[\emptyset]}\} \\ \stackrel{4.14(a)}{=} \{(\alpha(x_1), \dots, \alpha(x_k), b_1, \dots, b_\ell) : \alpha \text{ in } \text{ENUM}([\emptyset], \{<^v\}_{v \in V'})\}.$$

Since we can check in time  $O(1)$  (using the Boolean variable **start-is-fit** in the data structure) if  $[\emptyset]$  is fit and Lemma 4.14(b), we enumerate the tuples with delay  $t_d = O(|\text{vars}(Q)|)$  without repetition (this is guaranteed by Lemma 4.14(c)).

This concludes the proof of Theorem 3.3(b).

**Remark 4.15.** *As a remark, we show now that if the data structure is lexicographically ordered. Then the enumeration algorithm  $\text{ENUM}([\emptyset], \{<^v\}_{v \in V'})$  where  $x_j <_s^v x_{j'}$  if and only if  $j < j'$  and  $x_j, x_{j'} \in \text{child}(v)$  for all  $v \in V$ , enumerates the tuples in lexicographical order.*

*Let  $(\bar{a}_1, \bar{b}) \in Q(D)$  and  $(\bar{a}_2, \bar{b}) \in Q(D)$  where  $\bar{b} = (b_1, \dots, b_\ell)$  be two tuples such that  $\bar{a}_1$  will be enumerated before  $\bar{a}_2$ . Let  $j \in [k]$  be the smallest index with  $(\bar{a}_1)_j \neq (\bar{a}_2)_j$ . By definition of  $\{<^v\}_{v \in V'}$  is the outermost loop, that continued with the iteration between  $\bar{a}_1$  and  $\bar{a}_2$  were enumerated in the for loop of the form  $\widehat{t} \in \mathcal{L}_{x_j}^i$ . Since these lists are lexicographically ordered, it follows that  $(\bar{a}_1)_j < (\bar{a}_2)_j$  and thus  $\bar{a}_1 < \bar{a}_2$ .*

### 4.7. Enumerating the difference

The aim of this section is to show how to enumerate the tuples that are additionally in the result set and the tuples that do not belong to the result set any more after the database received some update steps, i.e., this chapter is devoted to prove Theorem 3.3 (c). For the remainder of this section we assume that  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$  is a  $q$ -hierarchical conjunctive query,  $\text{vars}(Q) = \{x_1, \dots, x_m\}$  with  $0 \leq k \leq m$ , and  $Q$  is of the form

$$Q = \{(x_1 \dots x_k, b_1, \dots, b_\ell) : \exists x_{k+1} \dots \exists x_m (\psi_1 \wedge \dots \wedge \psi_d)\}, \quad (4.5)$$

where  $b_1, \dots, b_\ell \in \mathbf{dom}$  and  $\psi_1, \dots, \psi_d$  are atomic queries of schema  $\sigma$ .

In the scenario here we need the system at one point to record the differences of the result set and after a couple of update steps we have the system start the enumeration of the difference. In this section we show how to enrich the data structure such that we can enumerate the difference with constant delay. Let  $D^-$  be the database at the start of the recording and let  $D^+$  be the database at the time we want to enumerate the difference, i.e.,  $D^+$  is the database obtained from  $D^-$  after the updates were processed. The aim is to enumerate the set  $Q(D^+) \setminus Q(D^-)$  (the “new” tuples in the result set)



#### 4.7. Enumerating the difference

and  $Q(D^-) \setminus Q(D^+)$  (the tuples that are not in the result any more). We call this routine the **diff**-routine.

The task of enumerating the difference has been considered in various papers in the sense of delta queries [59, 66, 67, 81]. A delta query is defined as a query whose result contains tuples that are additionally in the result and do not belong in the result of  $Q$  on  $D$  any more after the update. In [59] an update step might involve insertions and deletions of multiple tuples to the database instead of single tuples as stated in this thesis. Such updates can be simulated in the following way. Before proceeding an update step, we need the system to record the difference, and then to insert or delete the tuples into/from the database one by one. Note that here we use set semantics instead of bag semantics as in [59]. As a warm up, let us consider the following example.

**Example 4.16.** *As an example, let us consider the query  $Q$  and the database  $D_0$  in Example 4.4 as  $D^-$ . We need the system that to record the difference and receive the updates delete  $F(2, 2, 4)$ , insert  $E(3, 4)$ , insert  $F(3, 1, 2)$ , insert  $G(3, 1, 2)$ . Let  $D^+$  be the database we obtain after applying the updates and we then aim to start the enumeration of the difference. See Figure 4.6 for an illustration of  $D^+$ . The items that become fit after the update are marked yellow and the item that loses the fit status are marked red.*

*Note that the reason why there are tuples in  $Q(D^+) \setminus Q(D^-)$  is that new tuples were inserted to relations. In particular, we have  $a_0 = 3$  and  $(a_1 = 4$  or  $(a_2 = 1$  and  $a_3 = 2))$  for all  $(a_0, a_1, a_2, a_3) \in Q(D^+) \setminus Q(D^-)$ . This first idea is to enumerate the tuples  $(a_0, a_1, a_2, a_3) \in Q(D^+)$  with  $a_0 = 3$  and  $a_1 = 4$  and afterwards the tuples with  $a_0 = 3$  and  $a_2 = 1$  and  $a_3 = 2$ . But then we would enumerate the tuple  $(3, 4, 1, 2)$  twice. To get rid of this problem, we have to take into account that we enumerate in the second round only tuples with  $a_1 \neq 4$ .*

*To enumerate  $Q(D^-) \setminus Q(D^+)$  we only have to deal with tuples we removed from the relation and enumerate the tuples  $(a_0, a_1, a_2, a_3) \in Q(D^-)$  with  $a_0 = a_2 = 2$  and  $a_3 = 4$ .*

As we see in Example 4.16 here, there are some open questions for enumerating the difference.

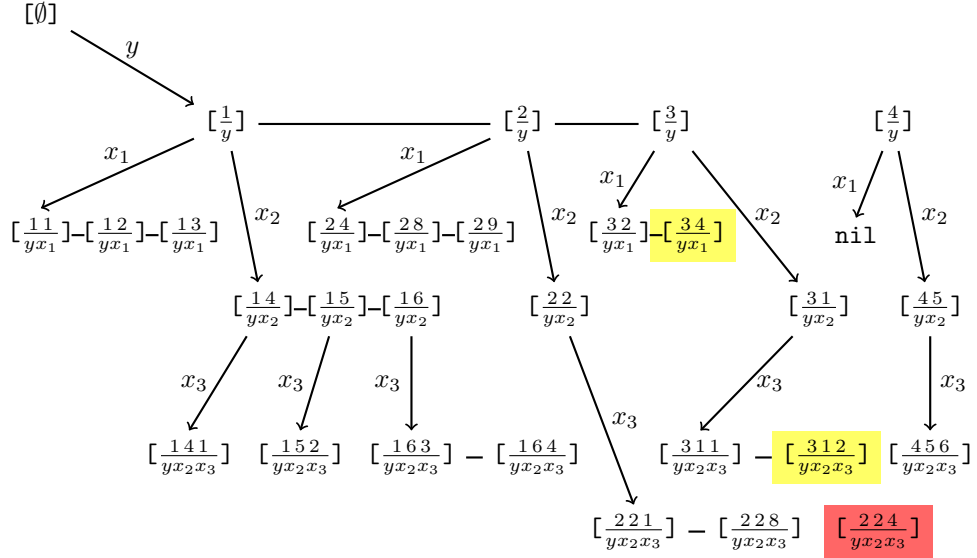
1. How to find the insertion (deletion) items that force new tuples being in  $Q(D^+) \setminus Q(D^-)$  or  $Q(D^-) \setminus Q(D^+)$ ?
2. How to ensure that we do not enumerate a tuple twice?

To answer question 1, note that it is not sufficient to simply iterate over all insertion steps and then iterate over the new tuples forced by the insertion. For instance, if we insert in our example database tuples  $(5, 1), \dots, (5, 100)$  to  $E$ , the difference  $Q(D^+) \setminus Q(D^-)$  is still empty but we can not enumerate this with constant delay. Thus, we use the fact, that whenever we have to take the insertion of a tuple to a relation into account for the difference, at least one item gets fit in the data structure.

To find these items efficiently we additionally store for every present item  $i$  in the data structure Boolean variables  $\text{fit}^+[i]$  and  $\text{fit}^-[i]$  and for every child  $u$  of  $v^i$  two

#### 4. Answering q-hierarchical conjunctive queries under updates

Figure 4.6.: Illustration from  $D^+$ .



lists in the data structure  $(\mathcal{L}_u^i)^+$  and  $(\mathcal{L}_u^i)^-$ . By default, the Boolean variables will be initialised to `false` and the lists will be initialised as empty lists. The aim is that the following for the Boolean variables and the lists holds.

$$\text{fit}^\circ[i] \text{ is true} \iff i \text{ is fit in } D^\circ \text{ and not in } D^\otimes \quad (4.6)$$

and

$$(\mathcal{L}_u^i)^\circ := \left\{ \iota : \begin{array}{l} \iota \in \mathcal{L}_u^i \text{ in } D^\circ \text{ and there is a } \beta \supseteq \alpha' \text{ such that } \mathbf{fit}^\circ[\beta] \text{ is true} \\ \text{and for all } w \in \mathbf{vpath}[v^{\lceil \beta \rceil}] \setminus \mathbf{vpath}[v^{\iota}] \text{ is } [\beta|_{\mathbf{vpath}[w]}] \text{ fit in } D^\circ \end{array} \right\} \quad (4.7)$$

for  $\circ \in \{+, -\}$  and  $\otimes \in \{+, -\} \setminus \{\circ\}$ . Intuitively, an item  $\iota$  belongs to a list  $(\mathcal{L}_u^i)^\circ$  if and only if there is a path from  $\iota$  to an item that is fit in  $D^\circ$  and not in  $D^\otimes$ . These lists will help us to find these items that are sufficient to enumerate the difference.

**Example 4.17.** Consider the setting of Example 4.16. Then  $\text{fit}^+[\frac{3}{y_{x_1}}] = \text{True}$  and for the other items  $i$  we have  $\text{fit}^+[i] = \text{False}$ .  $\text{fit}^+[\frac{3}{y_{x_2 x_3}}] = \text{True}$  and for the other items  $i$  is  $\text{fit}^+[i] = \text{False}$ . Furthermore,  $(\mathcal{L}_y^{\emptyset})^+ = \left\{ \frac{3}{y} \right\}$ ,  $(\mathcal{L}_{x_1}^{\frac{3}{y}})^+ = \left\{ \frac{3}{y_{x_1}} \right\}$ ,  $(\mathcal{L}_{x_2}^{\frac{3}{y}})^+ = \left\{ \frac{3}{y_{x_2}} \right\}$  and  $(\mathcal{L}_{x_3}^{\frac{3}{y_{x_2}}})^+ = \left\{ \frac{3}{y_{x_2 x_3}} \right\}$ . For all the other items  $\iota$  and for all  $u \in \text{child}(v^\iota)$  it is  $(\mathcal{L}_u^\iota)^+ = \emptyset$ . Moreover,  $(\mathcal{L}_y^{\emptyset})^- = \left\{ \frac{2}{y} \right\}$ ,  $(\mathcal{L}_{x_2}^{\frac{2}{y}})^- = \left\{ \frac{2}{y_{x_2}} \right\}$  and  $(\mathcal{L}_{x_3}^{\frac{2}{y_{x_2}}})^- = \left\{ \frac{2}{y_{x_2 x_3}} \right\}$ . For all the other items  $\iota$  and for all  $u \in \text{child}(v^\iota)$  is  $(\mathcal{L}_u^\iota)^- = \emptyset$ . As we see here, when following

#### 4.7. Enumerating the difference

the  $(\mathcal{L}_u^i)^\circ$  sets, we obtain a path from  $(\mathcal{L}_y^{\emptyset})^\circ$  to the items  $i$  where  $\text{fit}^\circ[i]$  is true for  $\circ \in \{+, -\}$ .

The list  $(\mathcal{L}_u^i)^-$  will be used to enumerate the tuples in  $Q(D^-) \setminus Q(D^+)$  since the fact that an item loses the fit-status implies that new tuples are in the result set. The idea is the following. An algorithm enumerates for  $\circ \in \{+, -\}$  the items  $i$  where  $\text{fit}^\circ[i]$  is true, using the  $(\mathcal{L}_u^i)^\circ$ -lists. Then we take for each such item  $i$  the assignment  $\alpha^i$  and enumerate the tuples  $(\beta(x_1), \dots, \beta(x_k), b_1, \dots, b_\ell)$  in  $Q(D^\circ)$  where  $\alpha^i \subseteq \beta$ , i.e., the output tuple coincides with the assignment of the item  $i$ . The next claim shows that these items help us to compute the tuples that belong to the difference set  $Q(D^\circ) \setminus Q(D^\otimes)$  where  $\otimes \in \{+, -\} \setminus \{\circ\}$ .

**Claim 4.18.** *Let  $Q$  be a  $q$ -hierarchical query and let  $D^-$  be a database and  $D^+$  be a database obtained from  $D^-$  after some updates on  $D^-$ . For every assignment  $\beta$  with  $(D^\circ, \beta) \models Q$  the following holds.*

$$\begin{aligned} &(\beta(x_1), \dots, \beta(x_k), b_1, \dots, b_\ell) \in Q(D^\circ) \setminus Q(D^\otimes) \iff \\ &\text{there is an item } i \text{ with } \alpha^i \subseteq \beta \text{ such that for all } w \in \text{free}(Q), \text{ is} \\ &\quad [\beta|_{\text{vpath}[w]}] \text{ fit in } D^\circ \text{ and } \text{fit}^\circ[i] \text{ is true} \end{aligned}$$

*Proof.*

$$\begin{aligned} &(\beta(x_1), \dots, \beta(x_k), b_1, \dots, b_\ell) \in Q(D^\circ) \setminus Q(D^\otimes) \\ \iff &(\beta(x_1), \dots, \beta(x_k), b_1, \dots, b_\ell) \in Q(D^\circ) \text{ and} \\ &(\beta(x_1), \dots, \beta(x_k), b_1, \dots, b_\ell) \notin Q(D^\otimes) \\ \iff &\text{for all } \psi \in \text{atoms}(Q) \text{ it holds } (D^\circ, \beta) \models \psi \text{ and} \\ &\text{there is a } \psi \in \text{atoms}(Q) \text{ such that } (D^\otimes, \beta) \not\models \psi \\ \iff &\text{for all } \psi \in \text{atoms}(Q) \text{ it holds } (D^\circ, \beta) \models \psi \text{ and} \\ &\text{there is a } \psi \in \text{atoms}(Q) \text{ such that } (D^\otimes, \beta|_{\text{vars}(\psi)}) \not\models \psi \\ \stackrel{\text{Def. fit}}{\iff} &\text{for all } w \in \text{free}(Q) \text{ is } [\beta|_{\text{vpath}[w]}] \text{ fit in } D^\circ \text{ and} \\ &[\beta|_{\text{vars}(\psi)}] \text{ is not fit in } D^\otimes \\ \iff &\text{for all } w \in \text{free}(Q) \text{ it holds } [\beta|_{\text{vpath}[w]}] \text{ is fit in } D^\circ \\ &\text{and there is an item } i \text{ in } D^\circ \text{ with } \alpha^i = \beta|_{\text{vars}(\psi)} \subseteq \beta \\ &\text{and } \text{fit}^\circ[i] \text{ is true} \end{aligned}$$

□

We will use the following claims to test efficiently if an item has to belong to a list  $(\mathcal{L}_u^i)^\circ$ .

**Claim 4.19.** *For all  $\circ \in \{+, -\}$ , for all present items  $\iota$  and their parent item  $i = \alpha^\iota|_{\text{vpath}[v^\iota]}$  the following holds:  $\iota \in (\mathcal{L}_u^i)^\circ$  if and only if  $\iota \in \mathcal{L}_u^i$  in  $D^\circ$  and  $\text{fit}^\circ[\iota]$  is true or there is a  $w \in \text{child}(v^\iota)$  such that  $(\mathcal{L}_w^\iota)^\circ \neq \emptyset$ .*

*Proof.* • For the “if”-direction let  $\iota \in (\mathcal{L}_u^i)^\circ$ . Clearly,  $\iota \in \mathcal{L}_u^i$  in  $D^\circ$ . Furthermore, there is a  $\beta \supseteq \alpha^\iota$  such that  $\text{fit}^\circ[\beta]$  is true and for all  $w \in \text{vpath}[v^{[\beta]}] \setminus \text{vpath}[v^\iota]$  is  $[\beta|_{\text{vpath}[w]}]$  fit in  $D^\circ$ .

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

- If  $\beta = \alpha^\iota$ , the claim holds since  $\text{fit}^\circ[\iota]$  is **true**.
- Otherwise, let  $w \in \text{child}(u)$  be the child such that  $w \in \text{dom}(\beta)$  and let  $\gamma := \beta|_{\text{vpath}[w]}$ . Since  $\gamma$  is fit in  $D^\circ$ , it holds that  $[\gamma] \in \mathcal{L}_w^\iota$  in  $D^\circ$  and there is a  $\beta \supseteq \gamma$  such that  $\text{fit}^\circ[\beta]$  is **true** and for all  $w \in \text{vpath}[v^{[\beta]}] \setminus \text{vpath}[\gamma]$  is  $[\beta|_{\text{vpath}[w]}]$  fit in  $D^\circ$ . Therefore,  $[\gamma] \in (\mathcal{L}_w^\iota)^\circ$  and, in particular,  $(\mathcal{L}_w^\iota)^\circ \neq \emptyset$ .
- For the “only if”-direction let us consider that  $\iota \in \mathcal{L}_u^i$  in  $D^\circ$ , and  $\text{fit}^\circ[\iota]$  is **true** or there is a  $w \in \text{child}(v^\iota)$  such that  $(\mathcal{L}_w^\iota)^\circ \neq \emptyset$ .
  - If  $\text{fit}^\circ[\iota]$  is **true**, the claim follows by the definition of  $(\mathcal{L}_u^i)^\circ$ .
  - Otherwise there is a  $w \in \text{child}(v^\iota)$  such that  $(\mathcal{L}_w^\iota)^\circ \neq \emptyset$ . Let  $\hat{\iota} \in (\mathcal{L}_w^\iota)^\circ$  arbitrary. Then, there is a  $\beta \supseteq \alpha^{\hat{\iota}}$  such that  $\text{fit}^\circ[\beta]$  is **true** for all  $w \in \text{vpath}[v^{[\beta]}] \setminus \text{vpath}[v^{\hat{\iota}}]$  we have that  $[\beta|_{\text{vpath}[w]}]$  is fit in  $D^\circ$ . Since  $\iota \in \mathcal{L}_u^i$  in  $D^\circ$ ,  $\iota$  is fit in  $D^\circ$  and we obtain  $\beta \supseteq \alpha^\iota$  and for all  $w \in \text{vpath}[v^{[\beta]}] \setminus \text{vpath}[v^\iota]$  we have that  $[\beta|_{\text{vpath}[w]}]$  is fit in  $D^\circ$ . Thus,  $\iota \in (\mathcal{L}_u^i)^\circ$ .

□

**Claim 4.20.** For all  $\otimes \in \{+, -\}$ , for all present items  $\iota$  and their parent item  $i = \alpha^\iota|_{\text{vpath}[v^\iota]}$  the following holds: It takes time  $O(|\text{child}(v^\iota)|)$  to test whether  $\iota \in (\mathcal{L}_u^i)^\otimes$ .

*Proof.* To test if  $\iota \in (\mathcal{L}_u^i)^\otimes$  is correct, it is sufficient by Claim 4.20, to test if  $\iota$  is fit in  $D^\circ$  and  $\text{fit}^\circ[\iota]$  is **true** or if there is a  $w \in \text{child}(v^\iota)$  such that  $(\mathcal{L}_w^\iota)^\circ$  is empty.

To test if  $\iota$  is fit in  $D^+$ , we simply test if it is fit in the current data structure. To test if  $\iota$  is fit in  $D^-$ , we test if  $\text{fit}^-[\iota]$  is **true** or  $\iota$  is fit in the current data structure and  $\text{fit}^+[\iota]$  is **false**. By definition of  $\text{fit}^\circ[\iota]$  we obtain the information if  $\iota$  is fit in  $D^-$ . Thus, the test can be done in  $O(|\text{vars}(Q)|)$ . □

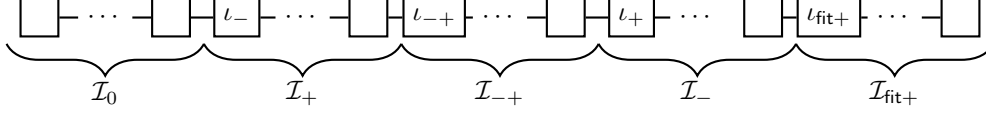
In the following Remark 4.21 it is described how to efficiently maintain the lists and to efficiently clean up the data structure, i.e., to obtain the data structure that we would have if we did not start the difference procedure.

**Remark 4.21.** Here we describe how to maintain the lists such that we can efficiently enumerate the items of the lists and how to clean up the data structure after the difference routine.

To avoid notational clutter, we define the following sets.

$$\begin{aligned}
 \mathcal{I}_0 &:= \mathcal{L}_u^i \setminus ((\mathcal{L}_u^i)^+ \cup (\mathcal{L}_u^i)^-) \\
 \mathcal{I}_- &:= (\mathcal{L}_u^i)^- \setminus (\{\iota \in (\mathcal{L}_u^i)^- : \text{fit}^-[\iota] \text{ is true}\} \cup (\mathcal{L}_u^i)^+) \\
 \mathcal{I}_{-+} &:= (\mathcal{L}_u^i)^- \cap (\mathcal{L}_u^i)^+ \\
 \mathcal{I}_+ &:= (\mathcal{L}_u^i)^+ \setminus (\{\iota \in (\mathcal{L}_u^i)^+ : \text{fit}^+[\iota] \text{ is true}\} \cup (\mathcal{L}_u^i)^-) \\
 \mathcal{I}_{\text{fit}^+} &:= \{\iota \in (\mathcal{L}_u^i)^+ : \text{fit}^+[\iota] \text{ is true}\} \\
 \mathcal{I}_{\text{fit}^-} &:= \{\iota \in (\mathcal{L}_u^i)^- : \text{fit}^-[\iota] \text{ is true}\}
 \end{aligned}$$

It is straightforward to verify that the sets above are pairwise disjoint. Note that every item in  $\iota \in \mathcal{I}_{-+}$  is fit in  $D^+$  and  $D^-$  (by Definition of  $(\mathcal{L}_u^i)^-$  and  $(\mathcal{L}_u^i)^+$ ) and thus by definition of  $\text{fit}^\circ[\iota]$  for  $\circ \in \{+, -\}$ , it follows that  $\text{fit}^+[\iota] = \text{fit}^-[\iota] = \text{false}$ .

Figure 4.7.: Illustration of  $\mathcal{L}_u^i$ .

Furthermore, it holds  $(\mathcal{L}_u^i)^+ \subseteq \mathcal{L}_u^i$  since every item that is fit in  $D^+$  is fit in the current database. Moreover,  $(\mathcal{L}_u^i)^- \setminus \mathcal{I}_{\text{fit}-} \subseteq \mathcal{L}_u^i$  since the items in  $(\mathcal{L}_u^i)^- \setminus \mathcal{I}_{\text{fit}-}$  are fit in  $D^+$  and in  $D^-$ . Using these facts it is straightforward to verify that

$$\mathcal{L}_u^i = \mathcal{I}_0 \cup \mathcal{I}_- \cup \mathcal{I}_{-+} \cup \mathcal{I}_+ \cup \mathcal{I}_{\text{fit}+}$$

We organize the list  $\mathcal{L}_u^i$  by the following way.

Let  $<$  be the linear order of the item in the list  $\mathcal{L}_u^i$ . We store the list  $\mathcal{L}_u^i$  as a concatenation of  $\mathcal{I}_0, \mathcal{I}_-, \mathcal{I}_{-+}, \mathcal{I}_+$  and  $\mathcal{I}_{\text{fit}+}$  (in this order). For all  $\circ \in \{+, -+, -, \text{fit}+\}$  let  $\iota_\circ$  be the first element of  $\mathcal{I}_\circ$  if it is not empty, i.e.,

$$\iota_\circ = \begin{cases} \text{the first item in } \mathcal{I}_\circ & \text{if } \mathcal{I}_\circ \neq \emptyset \\ \text{nil} & \text{if } \mathcal{I}_\circ = \emptyset \end{cases} \quad (4.8)$$

We call these items separator items.

In other words, if  $\iota_\circ \neq \text{nil}$  for all  $\circ \in \{+, -+, -, \text{fit}+\}$ , the items in  $\mathcal{L}_u^i$  are organized by the following way.

- The items in  $\iota \in \mathcal{L}_u^i$  with  $\iota < \iota_-$  belong to  $\mathcal{I}_0$ .
- The items in  $\iota \in \mathcal{L}_u^i$  with  $\iota_- \leq \iota < \iota_{-+}$  belong to  $\mathcal{I}_-$ .
- The items in  $\iota \in \mathcal{L}_u^i$  with  $\iota_{-+} \leq \iota < \iota_+$  belong to  $\mathcal{I}_{-+}$ .
- The items in  $\iota \in \mathcal{L}_u^i$  with  $\iota_+ \leq \iota < \iota_{\text{fit}+}$  belong to  $\mathcal{I}_+$ .
- The items in  $\iota \in \mathcal{L}_u^i$  with  $\iota_{\text{fit}+} \leq \iota$  belong to  $\mathcal{I}_{\text{fit}+}$ .

See Figure 4.7 for an illustration.

To maintain the order above, we use Algorithm 4 whenever we insert an item  $\iota$  to  $\mathcal{L}_u^i$  and Algorithm 5 whenever we delete an item  $\iota$  from  $\mathcal{L}_u^i$ . If an item modifies the membership to one of the sets  $\mathcal{I}_\circ$ , we use Algorithm 5 to delete  $\iota$  from  $\mathcal{L}_u^i$  and afterwards Algorithm 4 to insert the item  $\iota$  to  $\mathcal{L}_u^i$ . In these algorithms we let  $- < -+ < + < \text{fit}+$ .

It is straightforward to verify that we have the order described above if an item is inserted to the list  $\mathcal{L}_u^i$ , deleted from the list or the membership from one of the  $\mathcal{I}_\circ$  is modified.

To test whether an item  $\iota$  belongs to a list  $(\mathcal{L}_u^i)^\circ$ , we use Claim 4.20. This can be used to get the element  $\circ \in \{0, -, -+, +, \text{fit}+, \text{fit}-\}$  for which  $\iota \in \mathcal{I}_\circ$ . Thus, line 1

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

---

**Algorithm 4** An item  $\iota$  will be inserted to  $\mathcal{L}_u^i$ .

---

- 1: Test for which  $\circ \in \{0, -, -+, +, \text{fit}+, \text{fit}-\}$  it holds  $\iota \in \mathcal{I}_\circ$ . Let  $\circ$  be the corresponding element.
  - 2: **if**  $\iota \in \mathcal{I}_{\text{fit}-}$  **then**
  - 3:     Insert  $\iota$  to the list  $\mathcal{I}_{\text{fit}-}$ .
  - 4: **if**  $\iota_\circ \neq \text{nil}$  **then**
  - 5:     Insert  $\iota$  as the successor of  $\iota_\circ$ .
  - 6: **else**
  - 7:      $\iota_\circ := \iota$
  - 8:     **if** there is a  $\oplus \in \{-, -+, +, \text{fit}+\}$  with  $\circ < \oplus$  and  $\iota_\oplus \neq \text{nil}$  **then**
  - 9:         Let  $\otimes = \min \{\oplus \in \{-, -+, +, \text{fit}+\} : \circ < \oplus, \iota_\oplus \neq \text{nil}\}$
  - 10:         Insert  $\iota$  as the previous element of  $\iota_\otimes$ , i.e., as the successor element of the previous of  $\iota_\otimes$  if it exists or as the first element in the list otherwise.
  - 11:     **else**
  - 12:         Insert  $\iota$  as the last element of  $\mathcal{L}_u^i$ .
- 

---

**Algorithm 5** An item  $\iota$  will be deleted from  $\mathcal{L}_u^i$ .

---

- 1: **if** there is a  $\circ \in \{-, -+, +, \text{fit}+\}$  such that  $\iota_\circ = \iota$  **then**
  - 2:     **if** the successor element is an separator item, i.e., an item  $\iota_\oplus$  with  $\oplus \in \{-+, +, \text{fit}+\}$  **then**
  - 3:          $\iota_\circ := \text{nil}$ .
  - 4:     **else**
  - 5:         Define  $\iota_\circ$  as the successor item of  $\iota$ .
  - 6: Delete  $\iota$  from  $\mathcal{L}_u^i$ .
-

#### 4.7. Enumerating the difference

can be proceed in time  $O(\text{poly}(Q))$ . The other lines in Algorithm 4 and all lines in Algorithm 5 take time  $O(1)$ . Thus, it takes time  $O(\text{poly}(Q))$  to insert or delete an element in the list.

In the remainder of this section it is necessary to efficiently enumerate the elements in  $(\mathcal{L}_u^i)^-$ ,  $(\mathcal{L}_u^i)^+$  and  $\mathcal{L}_u^i \setminus \mathcal{I}_{\text{fit}+}$ . This can be done with the separator items:

- To enumerate  $(\mathcal{L}_u^i)^-$ , start with  $\iota_-$  (or  $\iota_{-+}$  if  $\iota_-$  is **nil**) and iterate through the elements until an item  $\iota_\circ$  with  $\circ \in \{+, \text{fit}+\}$  is reached. If  $\iota_- = \iota_{-+} = \text{nil}$ , then  $(\mathcal{L}_u^i)^- \setminus \mathcal{I}_{\text{fit}-} = \emptyset$ . Afterwards, we enumerate the set  $\mathcal{I}_{\text{fit}-}$ .
- To enumerate  $(\mathcal{L}_u^i)^+$ , start with  $\iota_{-+}$  (or  $\iota_+$  if  $\iota_{-+}$  is **nil**) and iterate through the elements until the end.
- To enumerate  $\mathcal{L}_u^i \setminus \mathcal{I}_{\text{fit}+}$ , start with the first item of  $\mathcal{L}_u^i$  and iterate until we reach  $\iota_{\text{fit}+}$  or, if  $\iota_{\text{fit}+}$  is **nil**, until the end of the list.

For every separator item  $\iota$ , we store a pointer to  $\iota$  together with a timestamp. The timestamp is equal to the number of times we request the enumeration of the difference.

The timestamp reveals if  $\iota$  is the current separator item. A separator item is not **nil** if and only if the timestamp equals the current timestamp. If we want to set the new separator item, we overwrite the timestamp. We use the same trick for the timestamps of  $\text{fit}^\circ[i]$ . Instead of  $\text{fit}^\circ[i] = \text{true}$  we store the last timestamp where we have set this value to **true**. Then, we know that  $\text{fit}^\circ[i]$  is set to **true** if and only if  $\text{fit}^\circ[i]$  has the current timestamp. This is convenient to clean up the data structure after enumerating the difference, i.e., to recover the data structure that we would obtain if we did not use the difference method. When resetting the difference procedure, we obtain that the timestamps get old and  $\text{fit}^\circ[i] = \text{false}$  for all items  $i$ . Furthermore, we have that the separator items are not the current separator items when the timestamps are getting old and the list  $\mathcal{I}_\circ$  is automatically empty for all  $\circ \in \{-, -+, +, \text{fit}+\}$ . For the set  $\mathcal{I}_{\text{fit}-}$  we also store a timestamp. The list is only a valid list if the stored timestamp is equal to the current timestamp. If we start a new list, we update the timestamp with a new list. If the timestamp is old, the list can be identified as empty. To clean up the data structure, we use the following idea. We store a list for every item  $\iota$  for that  $\text{fit}^-[\iota]$  is **true**. Whenever an item gets this status, we store the item at the end of the list and we remove the first item if the timestamp is old. Note that since we insert the items at the end of the list, it is sorted with respect to the timestamps.

We say that the data structure is correct if the conditions for the Boolean variables in Equation 4.6 and the lists in Equation 4.7 holds.

To ensure that the data structure is correct we execute, whenever an item gets fit during an update step, Algorithm 6 with  $\circ := +$  and, whenever an item loses the fit status during an update step, Algorithm 6 with  $\circ := -$ .

To show that Algorithm 6 is correct, we have to show that the algorithm ensures the correctness of the data structure.

**Claim 4.22.** 1. The data structure is correct if  $D^+ = D^-$ .

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

---

**Algorithm 6** The fit status of  $[\frac{a_1 \dots a_m}{v_1 \dots v_m}]$  changed.

---

```

1: Input:  $\circ \in \{+, -\}$ , item  $[\frac{a_1 \dots a_m}{v_1 \dots v_m}]$  whose fit status was changed.
2: Let  $\otimes \in \{+, -\} \setminus \{\circ\}$ 
3: Let  $\iota_j := [\frac{a_1 \dots a_j}{v_1 \dots v_j}]$  for all  $j \in \{1, \dots, m\}$ .
4: if  $\text{fit}^\otimes[i_m]$  is true then
5:    $\text{fit}^\otimes[i_m] \leftarrow \text{false}$ 
6:   for  $j = m$  to 1 do
7:     if  $i_j \in (\mathcal{L}_u^{i_{j-1}})^\otimes$  but it is not correct then
8:       Remove  $i_j$  from  $(\mathcal{L}_u^{i_{j-1}})^\otimes$ 
9: else if  $\text{fit}^\circ[i_m]$  is false then
10:   $\text{fit}^\circ[i_m] \leftarrow \text{true}$ 
11:  for  $j = m$  to 1 do
12:    if  $i_j$  is fit in  $D^\circ$  then insert  $i_j$  to  $(\mathcal{L}_{v_{j-1}}^{i_{j-1}})^\circ$ 
13:    else break ▷ Stop the for-loop

```

---

2. If the data structure is correct and an item gets fit (loses the fit status), Algorithm 6 with  $\circ := +$  (with  $\circ := -$ , resp.) ensures that the data structure is still correct.

*Proof.* Clearly, the claim holds for  $D^+ = D^-$  since we initialise  $(\mathcal{L}_u^i)^+$  and  $(\mathcal{L}_u^i)^-$  as empty lists and  $\text{fit}^+[i]$  and  $\text{fit}^-[i]$  as **false**.

For the second part we consider that the data structure is correct and ...

1. ... an item  $i_m = [\frac{a_1 \dots a_m}{v_1 \dots v_m}]$  gets fit. For all  $j \in [m]$  let  $i_j$  be the item  $[\frac{a_1 \dots a_j}{v_1 \dots v_j}]$ .
  - *Case 1a):* The item  $i_m$  was fit in  $D^-$ . Before  $i_m$  gets fit, it was fit in  $D^-$  and not in  $D^+$  and, in particular,  $\text{fit}^-[i_m]$  is **true** and has to be changed to **false**. Because of this fact, we have to check for all  $j \in [m]$  if  $i_j$  is still allowed to be in  $(\mathcal{L}_{v_j}^{i_{j-1}})^-$ . If not, we remove the item from the list. This is done by Algorithm 6 using  $\circ := +$ .
  - *Case 1b):* The item  $i_m$  was not fit in  $D^-$ . Before  $i_m$  gets fit, it was not fit in  $D^-$  and not in  $D^+$  and, in particular,  $\text{fit}^+[i_m]$  is **false** and has to be changed to **true**. Because of this fact, we have to add fit items  $i_j$  to  $(\mathcal{L}_{v_j}^{i_{j-1}})^+$  if there is a path over the  $+$ -lists to  $i_m$ . This is done by Algorithm 6 using  $\circ = +$ .
2. ... an item  $i_m = [\frac{a_1 \dots a_m}{v_1 \dots v_m}]$  loses the fit status. For all  $j \in [m]$  let  $i_j$  be the item  $[\frac{a_1 \dots a_j}{v_1 \dots v_j}]$ .
  - *Case 2a):* The item  $i_m$  was fit in  $D^-$ . Before  $i_m$  loses the fit status, it was fit in  $D^-$  and in  $D^+$  and, in particular,  $\text{fit}^-[i_m]$  is **false** and has to be changed to **true**. Because of this fact, we have to add fit items  $i_j$  in  $D^-$  to  $(\mathcal{L}_{v_j}^{i_{j-1}})^-$  if there is a path over the  $-$ -lists to  $i_m$ . This is done by Algorithm 6 using  $\circ = -$ .



#### 4.7. Enumerating the difference

- *Case 2b*): The item  $i_m$  was not fit in  $D^-$ . Before  $i_m$  loses the fit status, it was not fit in  $D^-$  but fit in  $D^+$  and, in particular,  $\text{fit}^+[i_m]$  is **true** and has to be changed to **false**. Because of this fact, we have to check for all  $j \in [m]$  whether  $i_j$  is still allowed to be in  $(\mathcal{L}_{v_j^{j-1}}^+)^+$ . If not, we remove the item from the list. This is done by Algorithm 6 using  $\circ = -$ .

Note that by the definition of the  $\mathcal{L}^\circ$  list it follows that if the fit status of an item  $i$  changes, it can only change the membership of an item  $\iota$  to the  $(\mathcal{L}_{v^\iota}^\circ)^\circ$  list for  $\circ \in \{+, -\}$  if  $\alpha^\iota \subseteq \alpha^i$ . Therefore, it is sufficient to consider only these items in the algorithm.  $\square$

Now, we show that updates can be efficiently done in constant time, i.e., Algorithm 6 takes constant time.

**Claim 4.23.** *Algorithm 6 takes time  $O(\text{vars}(Q)^2)$*

*Proof.* Since  $(\mathcal{L}_u^i)^-$  and  $(\mathcal{L}_u^i)^+$  are lists, lookup, insert and remove from the list can be done in  $O(|\text{vars}(Q)|)$  and changing the Boolean variables  $\text{fit}^+[i]$  and  $\text{fit}^-[i]$  can also be done in  $O(1)$ . When applying Claim 4.20, it takes time  $O(|\text{vars}(Q)|)$ . Note that the information whether there is a  $w \in \text{child}(v^\iota)$  such that  $(\mathcal{L}_w^\iota)^\circ$  is empty is correct since the for-loop goes from  $d$  to 1. Thus, the test can be done in  $O(|\text{vars}(Q)|)$ . All in all, Algorithm 6 takes time  $O(|\text{vars}(Q)|^2)$ .  $\square$

We now show how to use them to enumerate the difference.

We say that an item  $i$  is  $\circ$ -redundant if and only if there is a  $\beta \subset \alpha^i$  such that  $\text{fit}^\circ[[\beta]]$  is **true** for all  $\circ \in \{+, -\}$ . The main idea is to enumerate all non- $\circ$ -redundant items  $i$  and every extension of  $\alpha^i$ . A  $\circ$ -redundant item is not relevant, since for example if we have an item  $i$  and an item  $\iota$  with  $\alpha^i \subset \alpha^\iota$ , then we enumerate every extension of  $\alpha^\iota$  when we enumerate every extension of  $\alpha^i$ . Algorithm 7 shows how to enumerate with constant delay relevant non- $\circ$ -redundant items.

---

**Algorithm 7** Enumerate non- $\circ$ -redundant items  $i$  with  $\text{fit}^\circ[i]$ .

---

```

1: procedure ENUMFIT( $\circ, i$ )
2:   Input: Present item  $i$  and  $\circ \in \{+, -\}$ .
3:   if  $\text{fit}^\circ[i]$  is true then yield  $i$ 
4:   else
5:     for  $u \in \text{child}(v^i)$  do
6:       for  $\iota \in (\mathcal{L}_u^i)^\circ$  do
7:         yield ENUMFIT( $\circ, \iota$ )
8:   output EOE

```

---

The following claim states that the algorithm is correct.

**Claim 4.24.** *For all present items  $i$  in the data structure it holds:  $[\beta]$  is yielded by EnumFit( $\circ, [\emptyset]$ ) in Algorithm 7 if and only if*

- $\text{fit}^\circ[[\beta]]$  is **true**,

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

- $[\beta]$  is not  $\circ$ -redundant and
- for all  $w \in \text{vpath}[v^{[\beta]}]$  we have  $[\beta|_{\text{vpath}[w]}]$  fit in  $D^\circ$ .

*Proof.*

$[\beta]$  is yielded by  $\text{EnumFit}(\circ, [\emptyset])$

$$\begin{aligned} & \text{there are items } i_0, \dots, i_m \text{ and } v_0, \dots, v_m \in \text{vars}(Q) \text{ with } v_j = v^{i_j} \\ \iff & \text{ and } i_j \in (\mathcal{L}_{v_{j-1}}^{i_{j-1}})^\circ \text{ for all } j \in [m] \text{ such that } \text{fit}^\circ[\beta] \text{ is true} \\ & \text{and } \text{fit}^\circ[i_j] \text{ is false for all } j \in [d-1] \\ \iff & \text{ for all } w \in \text{vpath}[v^{[\beta]}] \text{ is } [\beta|_{\text{vpath}[w]}] \text{ fit in } D^\circ \text{ and} \\ & \text{fit}^\circ[\beta] \text{ is true and } [\beta] \text{ is not } \circ\text{-redundant} \end{aligned}$$

The first equation follows by the construction of the algorithm and the second equivalence from the definition of  $(\mathcal{L}_u^i)^\circ$  and  $\circ$ -redundant.  $\square$

The following claim states that the items will be enumerated with constant delay.

**Claim 4.25.**  $\text{EnumFit}(\circ, [\emptyset])$  enumerates the items with delay  $O(|\text{vars}(Q)|^2)$ .

*Proof.* By Claim 4.19 we obtain that on input of an item  $i$  that belongs to a set  $(\mathcal{L}_u^i)^\circ$  we have that  $\text{fit}^\circ[i]$  is true or there is a  $w \in \text{child}(v^i)$  such that  $(\mathcal{L}_w^i)^\circ \neq \emptyset$ . In particular, the algorithm will, upon input of such an item, yield at least one item. Since the recursion depth is bounded by  $|\text{vars}(Q)|$  the delay is  $O(|\text{vars}(Q)|^2)$ . If there is for the start-item  $[\emptyset]$  no  $w \in \text{child}(v_{\text{root}})$  with  $(\mathcal{L}_w^{[\emptyset]})^\circ \neq \emptyset$ , the algorithm takes time  $|\text{child}(v_{\text{root}})|$ . All the other items that  $\text{EnumFit}$  receives as input, belong to a list  $(\mathcal{L}_u^i)^\circ$ .  $\square$

To get rid of the problem that we might enumerate a tuple twice, the following definition will be convenient:

For all assignments  $\alpha$  and all  $\circ \in \{+, -\}$  let

$$D_{\text{VAR}}^\circ(\alpha) := \min \{v^i : \text{fit}^\circ[i] \text{ is true and } i \text{ is not redundant and } \alpha^i \subseteq \alpha\}$$

if it exists and  $D_{\text{VAR}}^\circ(\alpha) := \max \text{vars}(Q)$  otherwise.

**Example 4.26.** Recall Example 4.16. Then  $D_{\text{VAR}}^+(\mathcal{L}_{yx_1}^{\frac{3}{4}}J) = x_1$ ,  $D_{\text{VAR}}^+(\mathcal{L}_{yx_2}^{\frac{3}{1}}J) = x_3$ ,  $D_{\text{VAR}}^+(\mathcal{L}_{yx_2x_3}^{\frac{3}{1} \frac{2}{2}}J) = x_3$  and  $D_{\text{VAR}}^+(\mathcal{L}_{yx_1x_2x_3}^{\frac{3}{4} \frac{1}{2} \frac{2}{2}}J) = x_1$ .

The algorithm for enumerating the difference is the following:

---

**Algorithm 8** Enumerate the difference.

---

```

1: function DIFFENUM( $i, \hat{\iota}, \circ$ )
2:   Let  $u_1 < \dots < u_\ell$  be the children of  $v^i$  that belong to  $\text{free}(Q)$ .
3:   if  $v^i$  is a leaf then
4:     yield  $\alpha^i$ 
5:   else
6:     for  $\iota_1 \in \text{ES}_{u_1}^\circ(i, \hat{\iota})$  do
7:       for  $\alpha_1 \in \text{DIFFENUM}(\iota_1, \hat{\iota}, \circ)$  do
8:          $\vdots$ 
9:       for  $\iota_\ell \in \text{ES}_{u_\ell}^\circ(i, \hat{\iota})$  do
10:        for  $\alpha_\ell \in \text{DIFFENUM}(\iota_\ell, \hat{\iota}, \circ)$  do
11:          yield  $\alpha^i \cup \alpha_1 \cup \dots \cup \alpha_\ell$ 
12:
13: for  $\hat{\iota} \in \text{ENUMFIT}(\circ, [\emptyset])$  do
14:   for  $\alpha \in \text{DIFFENUM}([\emptyset], \hat{\iota}, \circ)$  do
15:     print  $(\alpha(x_1), \dots, \alpha(x_\ell), b_1, \dots, b_\ell)$ 

```

---

where

$$\text{ES}_v^\circ(i, \hat{\iota}) := \begin{cases} \{[\alpha|_{\text{vpath}[v]}]\} & \text{if } v \in \text{path}[v^{\hat{\iota}}], \\ \mathcal{L}_v^i \cup \{\iota \in (\mathcal{L}_v^i)^\circ : \text{fit}^\circ[\iota] = \text{true}\} & \text{if } v \notin \text{path}[v^{\hat{\iota}}] \text{ and } v > D_{\text{VAR}}^\circ(\alpha^i), \\ \mathcal{L}_v^i \setminus \{\iota \in (\mathcal{L}_v^i)^\circ : \text{fit}^\circ[\iota] = \text{true}\} & \text{otherwise.} \end{cases}$$

The set  $\text{ES}_v^\circ(i, \hat{\iota})$  will help us to avoid that a tuple will be enumerated twice.

**Example 4.27.** Recall Example 4.16 and let  $\iota_1 := [\frac{3}{y x_1} \frac{4}{}]$  and  $\iota_2 := [\frac{3}{y x_2 x_3} \frac{1}{2} \frac{2}{}]$ . The sets  $\text{ES}_v^+(i, \iota_j)$  enumerate tuples which coincide only with the assignments  $\alpha^{\iota_j}$  for  $j \in \{1, 2\}$ . Since  $\text{ES}_y^+(\emptyset, \iota_1) = \{[\frac{3}{y}]\}$ , the enumeration will be forced to use the assignment  $[\frac{3}{y}]$  for the enumeration. Then,  $\text{ES}_{x_1}^+([\frac{3}{y}], \iota_1) = \{[\frac{3}{y x_1} \frac{4}{}]\}$ , and  $\text{ES}_{x_2}^+([\frac{3}{y}], \iota_1) = \mathcal{L}_{x_2}^{[\frac{3}{y}]}$  since  $x_2 > x_1 = D_{\text{VAR}}^+(\frac{3}{y x_1})$ . This ensures that all tuples were enumerated that coincide with  $\alpha^{\iota_1}$ . Let us now consider  $\iota_2$ . Then,  $\text{ES}_y^+(\emptyset, \iota_2) = \{[\frac{3}{y}]\}$ ,  $\text{ES}_{x_1}^+([\frac{3}{y}], \iota_2) = \mathcal{L}_{x_1}^{[\frac{3}{y}]} \setminus \left\{ \iota \in (\mathcal{L}_{x_2}^{[\frac{3}{y}]})^\circ : \text{fit}^\circ[\iota] = \text{true} \right\} = \{[\frac{3}{y x_1} \frac{2}{}]\}$ , since  $\text{fit}^+[\frac{3}{y x_1} \frac{4}{}]$  is true. This ensures that the tuple  $(3, 4, 1, 2)$  will not be enumerated twice.

The next set will help us to describe the set of assignments that will be yielded by the DIFFENUM procedure. For all present items  $i$  and all assignments  $\gamma$  and all  $\circ \in \{+, -\}$  let:

$$(\Delta_\gamma^i)^\circ := \left\{ \beta \supseteq \alpha^i : \begin{array}{l} \text{dom}(\beta) = \bigcup_{\psi \in \text{atoms}(v^i)} \text{vars}(\psi), \\ (D^\circ, \beta) \models \psi \text{ for all } \psi \in \text{atoms}(v^i), \\ \text{for all } v \in \text{dom}(\beta) \cap \text{dom}(\gamma) \text{ is } \beta(v) = \gamma(v), \\ D_{\text{VAR}}^\circ(\beta) \geq v^{[\gamma]} \end{array} \right\}$$

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

Here,  $\gamma$  is the assignment we take into account for the enumeration, i.e., the aim is to enumerate all assignments  $\beta$  with  $\gamma \subseteq \beta$ . The set  $(\Delta_\gamma^i)^\circ$  is the set of assignments that are extensions of  $\alpha^i$  that we consider, when enumerating with  $\gamma$ .

**Example 4.28.** Let us consider the assignments of the two items that gets fits after the updates in Example 4.16, i.e.,  $\gamma_1 := \frac{3\ 4}{yx_1}$  and  $\gamma_2 := \frac{3\ 1\ 2}{yx_2x_3}$  and let  $i = [\frac{3}{y}]$ . Then,  $\frac{3\ 4\ 1\ 2}{yx_1x_2x_3} \in (\Delta_{\gamma_1}^i)^\circ$  and  $\frac{3\ 4\ 1\ 2}{yx_1x_2x_3} \notin (\Delta_{\gamma_2}^i)^\circ$  since  $D_{VAR}^+(\lceil \frac{3\ 4\ 1\ 2}{yx_1x_2x_3} \rceil) = x_1 < x_3 = D_{VAR}^+(\gamma_2)$ . Thus, when we first enumerate the tuples that coincide with  $\gamma_1$ , the tuple  $(3, 4, 1, 2)$  appears and when we enumerate the tuples that coincides with  $\gamma_2$ , we do not enumerate  $(3, 4, 1, 2)$  for a second time.

We now show that this is the set of assignments yielded by the DIFFENUM procedure:

**Claim 4.29.** For all assignments  $\gamma$  with  $\text{dom}(\gamma) = \text{vpath}[v]$  for a  $v \in \text{vars}(Q)$  and for all items  $i$ , where

- for all  $v \in \text{dom}(\gamma) \cap \text{dom}(\alpha^i)$  we have  $\gamma(v) = \alpha^i(v)$ ,
- $D_{VAR}^\circ(\alpha^i) > v^{\lceil \gamma \rceil}$  and
- $i$  is fit in  $D^\circ$ ,

the following holds for all  $\circ \in \{+, -\}$ .  $\text{DIFFENUM}(i, \lceil \gamma \rceil, \circ)$  enumerates the set  $(\Delta_\gamma^i)^\circ$ .

*Proof.* Let  $T_Q$  be the  $q$ -tree of  $Q$  from which the data structure was constructed from and let  $T_{\text{free}}$  be the induced subgraph of  $T_Q$  induced on  $\text{free}(Q)$ . We show this claim by induction over the height of  $i$  in  $T_{\text{free}}$ . For the induction base, let us consider an item  $i$  where  $v^i$  is a leaf in the  $q$ -tree of  $Q$ . By assumption and since  $i$  is fit in  $D^\circ$  it holds  $(\Delta_\gamma^i)^\circ = \{\alpha^i\}$ .

For the inductive step, let us consider an item of height  $> 0$ . Let  $u_1, \dots, u_s$  denote

#### 4.7. Enumerating the difference

the children of  $v^i$  in  $T_{\text{free}}$ . Then, for all assignments  $\beta$  the following holds.

$$\begin{aligned}
\beta \in (\Delta_\gamma^i)^\circ &\stackrel{\text{Def. } (\Delta_\gamma^i)^\circ}{\iff} \begin{aligned} &\beta \supseteq \alpha^i, \text{ dom}(\beta) = \bigcup_{\psi \in \text{atoms}(v^i)} \text{vars}(\psi), \\ &(D^\circ, \beta) \models \psi \text{ for all } \psi \in \text{atoms}(v^i), \\ &\text{for all } v \in \text{dom}(\beta) \cap \text{dom}(\gamma) \text{ is } \beta(v) = \gamma(v), \\ &\text{and } D_{\text{VAR}}^\circ(\beta) \geq v^{[\gamma]} \end{aligned} \\
&\stackrel{(*)}{\iff} \begin{aligned} &\beta \supseteq \alpha^i, \text{ dom}(\beta) = \bigcup_{\psi \in \text{atoms}(v^i)} \text{vars}(\psi), \\ &(D^\circ, \beta) \models \psi \text{ for all } \psi \in \text{exatoms}(v^i), \text{ and} \\ &\text{for all } j \in [s] \text{ holds } (D^\circ, \beta|_{\text{vpath}[v] \cup \text{succ}(u_j)}) \models \psi \\ &\quad \text{for all } \psi \in \text{atoms}(u_j), \\ &\text{for all } v \in \text{dom}(\beta) \cap \text{dom}(\gamma) \text{ is } \beta(v) = \gamma(v), \\ &\text{and } D_{\text{VAR}}^\circ(\beta) \geq v^{[\gamma]} \end{aligned} \\
&\iff \begin{aligned} &\beta \supseteq \alpha^i, \text{ dom}(\beta) = \bigcup_{\psi \in \text{atoms}(v^i)} \text{vars}(\psi), i \text{ is fit in } D^\circ \\ &\text{and for all } j \in [s] \text{ there is a } \iota \in \text{ES}_{u_j}^\circ(i, [\gamma]) \\ &\quad \text{such that } \beta|_{\text{vpath}[v] \cup \text{succ}(u_j)} \in (\Delta_\gamma^\iota)^\circ \text{ and} \\ &\text{for all } v \in \text{dom}(\beta) \cap \text{dom}(\gamma) \text{ is } \beta(v) = \gamma(v), \text{ and } D_{\text{VAR}}^\circ(\beta) \geq v^{[\gamma]} \end{aligned}
\end{aligned}$$

The equivalence  $(*)$  follows from the fact that

$$\text{atoms}(v^i) = \text{exatoms}(v^i) \cup \text{atoms}(u_1) \cup \dots \cup \text{atoms}(u_s).$$

By the induction hypothesis it follows that  $\beta|_{\text{vpath}[v] \cup \text{succ}(u_j)}$  will be enumerated by  $\text{DIFFENUM}(\iota, [\gamma], \circ)$ . It remains to show the last equivalence, i.e., for all  $u \in \text{child}(v^i)$  holds  $(D^\circ, \beta_u) \models \psi$  for all  $\psi \in \text{atoms}(u)$  if and only if for all  $u \in \text{child}(v^i)$  there is a  $\iota \in \text{ES}_u^\circ(i, [\gamma])$  such that  $\beta_u \in (\Delta_\gamma^\iota)^\circ$  where  $\beta_u := \beta|_{\text{vpath}[u] \cup \text{succ}(u)}$ .

For the “if” direction let  $u \in \text{child}(v)$  be arbitrary and assume that  $(D^\circ, \beta_u) \models \psi$  for all  $\psi \in \text{atoms}(u)$ . In the first step we show  $[\beta|_{\text{vpath}[u]}] \in \text{ES}_u^\circ(i, [\gamma])$ :

- *Case 1:*  $u \in \text{path}[v^{[\gamma]}]$ : Then  $\text{ES}_u^\circ(i, [\gamma]) = \{[\gamma|_{\text{vpath}[u]}]\}$ . Since  $\beta(v) = \gamma(v)$  for all  $v \in \text{dom}(\beta) \cap \text{dom}(\gamma)$  it holds that  $\gamma|_{\text{vpath}[u]} = \beta|_{\text{vpath}[u]}$  and therefore  $[\beta|_{\text{vpath}[u]}]$  is in  $\text{ES}_u^\circ(i, [\gamma])$ .
- *Case 2:*  $u \notin \text{path}[v^{[\gamma]}]$ : Since  $(D^\circ, \beta_u) \models \psi$  for all  $\psi \in \text{atoms}(u)$  it follows that  $\beta_u$  is fit in  $D^\circ$  and, in particular,  $[\beta|_{\text{vpath}[u]}]$  belongs to  $\mathcal{L}_u^{[\beta|_{\text{vpath}[u]}]}$  or  $\{\iota \in (\mathcal{L}_u^{[\beta|_{\text{vpath}[u]}]})^\circ : \text{fit}^\circ[\beta_u] \text{ is true}\}$ . Hence,  $[\beta|_{\text{vpath}[u]}] \in \text{ES}_u^\circ(i, [\gamma])$  if  $v \notin \text{path}[v^i]$  and  $v > D_{\text{VAR}}^\circ(\alpha^i)$ . Let us now consider the case that  $v < D_{\text{VAR}}^\circ(\alpha^i)$ . Assume for a contradiction that  $u < D_{\text{VAR}}^\circ(\gamma)$  and  $\text{fit}^\circ[\beta|_{\text{vpath}[u]}]$  is true. Then it follows that  $D_{\text{VAR}}^\circ(\beta_u) \leq u$  and thus  $D_{\text{VAR}}^\circ(\beta) \leq u < v^{[\gamma]}$  (since  $\beta|_{\text{vpath}[v]} \subseteq \beta$ ) and thus  $\beta \notin (\Delta_\gamma^i)^\circ$  which violates  $\beta \in (\Delta_\gamma^i)^\circ$ . Hence,  $[\beta|_{\text{vpath}[u]}] \in \text{ES}_u^\circ(i, [\gamma])$  if  $v \notin \text{path}[v^i]$  and  $v < D_{\text{VAR}}^\circ(\alpha^i)$ .

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

Let  $\iota := [\beta|_{\text{vpath}[u]}]$ . It remains to show that  $\beta_u \in (\Delta_\gamma^\iota)^\circ$ . It is known that  $(D^\circ, \beta_u) \models \psi$  for all  $\psi \in \text{atoms}(v^i)$ . Since for all  $v \in \text{dom}(\beta) \cap \text{dom}(\gamma)$  we have  $\beta(v) = \gamma(v)$  and because  $\beta_u \subseteq \beta$  for all  $v \in \text{dom}(\beta_u) \cap \text{dom}(\gamma)$  we have  $\beta_u(v) = \gamma(v)$ .

Assume for a contradiction that  $D_{\text{VAR}}^\circ(\beta_u) < v^{[\gamma]}$ . Then, there is an item  $i$  with  $\alpha^i \subseteq \beta_u$  with  $v^i = D_{\text{VAR}}^\circ(\beta_u)$  such that  $\text{fit}^\circ[i]$  is **true** and  $v^i < v^{[\gamma]}$ . Since  $\beta \supseteq \beta_u$ , also  $\alpha^i \subseteq \beta$  and therefore  $D_{\text{VAR}}^\circ(\beta) < v^{[\gamma]}$ . Thus, it follows that  $\beta \notin (\Delta_\gamma^i)^\circ$  which violates  $\beta \in (\Delta_\gamma^i)^\circ$ .

Therefore  $D_{\text{VAR}}^\circ(\beta_u) \geq v^{[\gamma]}$  and altogether  $\beta_u \in (\Delta_\gamma^\iota)^\circ$ .

For the “only if” direction, let us assume that  $i$  is fit and for all  $u \in \text{child}(v^i)$  there is an  $\iota \in \text{Es}_u^\circ$  such that  $\beta_u \in (\Delta_\gamma^\iota)^\circ$ . By the definition of  $(\Delta_\gamma^\iota)^\circ$  and since  $i$  is fit in  $D^\circ$ , it follows that  $(D^\circ, \beta) \models \psi$  for all  $\psi \in \bigcup_{u \in \text{child}(v)} \text{atoms}(u) \cup \text{exatoms}(v) = \text{atoms}(v)$ .  $\square$

Algorithm 8 does not enumerate a tuple twice since for all  $\iota, \iota'$  that will be yielded by Algorithm 7, it holds  $\alpha^\iota \not\subseteq \alpha^{\iota'}$  and  $\alpha^{\iota'} \not\subseteq \alpha^\iota$  and therefore  $(\Delta_\gamma^\iota)^\circ \cap (\Delta_\gamma^{\iota'})^\circ = \emptyset$ .

Recall, that we have constructed the data structure such that we can efficiently enumerate the set  $\mathcal{L}_u^i \setminus \{\iota : \text{fit}^+[i] \text{ is true}\}$  and the fact, that the delay takes time  $O(\text{vars}(Q)^2)$  can be shown analogously to the delay time of Algorithm 3.

The next claim shows that Algorithm 8 is correct.

**Claim 4.30.** *The tuples printed in line Algorithm 8 are exactly the tuples that belong to  $D^\circ \setminus D^\otimes$ .*

#### 4.8. Outputting the next smallest tuple

*Proof.*

$\bar{a}$  will be printed in Algorithm 8

$\Leftrightarrow^{(*)}$  there is a  $\hat{t} \in \text{EnumFit}(\circ, [\emptyset])$  and a  $\alpha \in \text{DiffEnum}([\emptyset, \hat{t}, \circ])$   
such that  $a = (\alpha(x_1), \dots, \alpha(x_\ell), b_1, \dots, b_\ell)$

$\Leftrightarrow^{(**)}$  there is an item  $\hat{t} = [\beta]$   
such that for all  $w \in \text{vpath}[v^{[\beta]}]$  is  $[\beta|_{\text{vpath}[w]}]$  fit in  $D^\circ$  and  
 $\text{fit}^\circ[\beta]$  is true and  $[\beta]$  is not redundant  
and a  $\alpha \in \left(\Delta_\beta^{[\emptyset]}\right)^\circ$  such that  $a = (\alpha(x_1), \dots, \alpha(x_\ell), b_1, \dots, b_\ell)$

$\Leftrightarrow^{(***)}$  there is an item  $\hat{t} = [\beta]$   
such that for all  $w \in \text{vpath}[v^{[\beta]}]$  is  $[\beta|_{\text{vpath}[w]}]$  fit in  $D^\circ$  and  
 $\text{fit}^\circ[\beta]$  is true and  $[\beta]$  is not redundant  
and a  $\alpha$  with  $\text{dom}(\alpha) = \text{vars}(Q)$  and  $(D^\circ, \alpha) \models Q$  and  
 $\beta \subseteq \alpha$  and  $D_{\text{VAR}}^\circ(\alpha) \geq v^{[\beta]}$  such that  $\bar{a} = (\alpha(x_1), \dots, \alpha(x_\ell), b_1, \dots, b_\ell)$

$\Leftrightarrow^{(****)}$  there is an assignment  $\alpha$  with  $\text{dom}(\alpha) = \text{vars}(Q)$  and  
 $\bar{a} = (\alpha(x_1), \dots, \alpha(x_\ell), b_1, \dots, b_\ell)$  and there is an item  $\hat{t} = [\beta]$   
with  $\beta \subseteq \alpha$  such that for all  $w \in \text{free}(Q)$  we have  
 $[\alpha|_{\text{vpath}[w]}]$  is fit in  $D^\circ$ , and  $\text{fit}^\circ[\hat{t}]$  is true

$\xLeftrightarrow{4.18} \bar{a} = (\alpha(x_1), \dots, \alpha(x_k), b_1, \dots, b_\ell) \in Q(D^\circ) \setminus Q(D^\otimes)$

Equation  $(*)$  follows from the construction of the algorithm, Equation  $(**)$  by Claim 4.24 and 4.29, Equation  $(***)$  by the Definition of  $\left(\Delta_\beta^{[\emptyset]}\right)^\circ$ . For the “if” direction in  $(****)$  Note that for all  $w \in \text{vars}(Q)$  we have  $[\alpha|_{\text{vpath}[w]}]$  is fit in  $D^\circ$  follows from  $(D^\circ, \alpha) \models Q$ . For the “only if” direction of  $(****)$  note that one can choose an item  $\hat{t}$  with  $D_{\text{VAR}}^\circ(\alpha) = v^{\hat{t}}$  and by Definition of  $D_{\text{VAR}}^\circ(\alpha)$  and the fact that the variables are sorted in pre-order, it follows that  $\hat{t}$  is not  $\circ$ -redundant.  $\square$

This concludes the proof of Theorem 3.3(c).

## 4.8. Outputting the next smallest tuple

In this section we will show that for a lexicographically ordered data structure, we can do the following routine. Upon input of a tuple  $\bar{a} \in \mathbf{dom}^{k+\ell}$ , output the following tuple

$$\max \{ \bar{c} : \bar{c} \leq_{\text{lex}} \bar{a}, \bar{c} \in Q(D) \}$$

if it exists, and **SmallerThanMinimum** otherwise.

The result is stated in the following lemma:

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

**Lemma 4.31.** *There is a dynamic algorithm that receives a  $q$ -hierarchical  $k$ -ary CQ  $Q$  and a  $\sigma$ -db  $D_0$  of size  $\|D_0\|$ , and computes within preprocessing time  $\text{poly}(Q) \cdot O(\|D_0\| \log(\|D_0\|))$  a data structure that can be updated in time  $\text{poly}(Q) \cdot O(\log(\|D\|))$  and allows the following:*

- (a) *enumerate the tuples in  $Q(D)$  in lexicographical order with delay  $\text{poly}(Q)$  and*
- (b) *upon input of a tuple  $\bar{b} \in \text{dom}^{k+\ell}$  output the tuple*

$$\max \{ \bar{a} \in Q(D) : \bar{a} \leq \bar{b} \}$$

*if it exists, or a **SmallerThanMinimum**-message otherwise in time  $t_l = \text{poly}(Q)$*

where  $D$  is the current database.

We have already seen that updating a lexicographically ordered data structure can be done in  $O(\log \|D\|)$  and it supports enumerating the query result in lexicographical order (see Remark 4.9 and 4.15). In particular, this data structure will be used for Lemma 4.31. Moreover, the routine in Lemma 4.31 (a) is also supported. The aim of this section is to show that Lemma 4.31 (b) holds. Before we discuss the general case, let us consider a running example.

**Example 4.32.** *Suppose, we have the lexicographically ordered data structure for  $Q$  on  $D_0$  from Example 4.4. Note that we obtain the data structure from Figure 4.3 where the items  $\lfloor \frac{2}{y x_2 x_3} \rfloor$  and  $\lfloor \frac{2}{y x_2 x_3} \rfloor$  are inverted in the  $x_3$ -list of  $\lfloor \frac{2}{y x_2 x_3} \rfloor$ .*

*Assume, we receive as input the tuple  $\bar{a} = (1, 2, 3, 1)$  and the aim is to compute the maximum tuple  $\bar{c}$  that belongs to the query result and is smaller than  $\bar{a}$ , i.e., the tuple  $\bar{c} = (1, 1, 6, 4)$ . Let  $\alpha$  be the assignment that coincides with  $\bar{a}$ , i.e.,  $\alpha = \frac{1}{y x_1 x_2 x_3} \frac{2}{y x_2 x_3} \frac{3}{y x_2 x_3} \frac{1}{y x_2 x_3}$ . In the first step, we compute two variables  $\tilde{v}$  and  $v$ . The variable  $\tilde{v} \in \{y, x_1, x_2, x_3\}$  is the minimal variable such that  $[\alpha|_{\text{vpath}[\tilde{v}]}]$  is not in the  $\tilde{v}$ -list of  $[\alpha|_{\text{vpath}[\tilde{v}]}]$ . In our example, this is  $\tilde{v} = x_2$ . Then, we let  $v \in \{y, x_1, x_2\} = \{w \in \text{vars}(Q) : w \leq x_2\}$  be the maximum variable such that there is an item  $\iota \in \mathcal{L}_v^{[\alpha|_{\text{vpath}[v]}]}$  with  $\alpha^\iota < \alpha(v)$ , and  $v$  is smaller than or equal to  $\tilde{v}$ . In our example,  $v = x_1$ , since there is an item  $\lfloor \frac{1}{y x_1} \rfloor$  that appears before  $\lfloor \frac{1}{y x_1} \rfloor$  in the  $x_1$ -list of  $\lfloor \frac{1}{y} \rfloor$ .*

*The idea is the following. Since the lists in the data structure are sorted, the enumeration algorithm in Algorithm 2 will enumerate the tuples in lexicographical order. Therefore, we ask for the maximum tuple, the enumeration algorithm will output, that is smaller than  $\bar{a}$ . Note that we have selected  $\tilde{v}$  as the minimum variable, such that  $[\alpha|_{\text{vpath}[\tilde{v}]}]$  is not in the  $\tilde{v}$ -list of  $[\alpha|_{\text{vpath}[\tilde{v}]}]$ . In particular, it follows that in the for-loop in line 4 we will not iterate  $\lfloor \frac{1}{y x_2} \rfloor$ . Since this is the minimum variable for which the condition holds, the for-loop in line 4 is the outermost for-loop where we do not find a corresponding item of the form  $[\alpha|_{\text{path}[u]}]$ . Therefore we upwardly try to find upwardly a for-loop that has a previous step for  $\lfloor \frac{1}{y x_2} \rfloor$ . This is not in line 4 the case since the first element in the  $x_2$ -list of  $\lfloor \frac{1}{y} \rfloor$  is  $\lfloor \frac{1}{y x_2} \rfloor$ . But the for-loop in line 3 has a previous element for  $\lfloor \frac{1}{y x_1} \rfloor$  in the  $x_1$ -list of  $\lfloor \frac{1}{y} \rfloor$ . This is exactly the condition that holds for  $v$ . Let us assume we set the for-loop in line 3 to the previous step, the set*



#### 4.8. Outputting the next smallest tuple

the for-loop in line 2 to the configuration that matches with  $\bar{a}$  and the for-loop in line 4 and 5 to the last configuration. Then, the tuple is smaller than  $\bar{a}$  and its successor tuple in the enumeration is lexicographically greater.

In other words, we construct the following assignment.

- For all  $v' < v$ , we set  $\gamma(v') = \alpha(v)$ . Then,  $\gamma(y) = 1$ .
- Then, we set  $\gamma(v)$  to the constant of the item that appears before  $\lceil \frac{1}{y x_1} \rceil$ , i.e.,  $\gamma(x_1) = 1$ .
- For all  $v' > v$ , we set  $\gamma(v') = a^\iota$ , where  $\iota$  is the last element of  $x_j$ -list of  $\lceil \gamma|_{\text{vpath}[x_j]} \rceil$ , i.e.,  $\gamma(x_2) = 6$  and  $\gamma(x_3) = 4$

We obtain  $\bar{c}$  by setting  $\bar{c} = (\gamma(y), \gamma(x_1), \gamma(x_2), \gamma(x_3)) = (1, 1, 6, 4)$ . If  $\bar{a}$  would appear in the enumeration algorithm, it would be followed by  $\bar{c}$ .

Now we consider the general case. In the reminder of this section we will show that Lemma 4.31 (b) holds.

For the remainder of this proof we assume that  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$  is a q-hierarchical conjunctive query,  $\text{vars}(Q) = \{x_1, \dots, x_m\}$  with  $0 \leq k \leq m$ , and  $Q$  is of the form

$$Q = \{(x_1 \dots x_k, b_1, \dots, b_\ell) : \exists x_{k+1} \dots \exists x_m (\psi_1 \wedge \dots \wedge \psi_d)\}, \quad (4.9)$$

where  $b_1, \dots, b_\ell \in \mathbf{dom}$  and  $\psi_1, \dots, \psi_d$  are atomic queries of schema  $\sigma$ . Let

$$\varphi := \exists x_{k+1} \dots \exists x_m (\psi_1 \wedge \dots \wedge \psi_d).$$

Let  $\bar{a}$  be the tuple we receive as input. The algorithm for computing the tuple

$$\bar{c} = \max \{\bar{d} \in Q(D) : \bar{d} \leq \bar{a}\}$$

is the following. If  $\bar{a}$  is smaller than the minimum element, we output the message **SmallerThanMinimum** and stop the algorithm. If  $\bar{a} \in Q(D)$  we output  $\bar{a}$  and stop the algorithm. Let  $\alpha$  be the assignment such that  $\alpha(x_j) := a_j$  for all  $j \in [k]$ . If  $\bar{a} \notin Q(D)$  and  $(D, \alpha) \not\models \varphi$  we construct an assignment  $\gamma$  such that for the number

$$p := \max \left\{ j \in [\hat{p}] : \text{there is an } \iota \in \mathcal{L}_{x_j}^{[\alpha|_{\text{vpath}[x_j]}]} \text{ such that } a^\iota < a_i \right\}$$

where

$$\hat{p} := \min \{ j \in [k] : [\alpha|_{\text{vpath}[x_j]}] \text{ is not in the } x_j\text{-list of } [\alpha|_{\text{vpath}[x_j]}] \}$$

the following holds:

- for all  $j < p$  is  $\gamma(x_j) = \alpha(x_j)$  and
- $\gamma(x_p) = \max \{ a^\iota : \iota \in \mathcal{L}_{x_\ell}^{[\alpha|_{\text{vpath}[x_\ell]}]} , a^\iota < a_i \}$  and

#### 4. Answering $q$ -hierarchical conjunctive queries under updates

- for all  $j > p$  is  $\gamma(x_j) = \alpha^\iota$  the maximum item  $\iota$  of the  $x_j$ -list of  $[\gamma|_{\text{vpath}[x_j]}]$ .

Note that an index  $p$  exists since  $\bar{a}$  is not smaller than the minimum. As a remark, note that  $x_p$  coincides with  $v$  in Example 4.32 and  $x_{\bar{p}}$  with  $\tilde{v}$ . Note that the elements can be computed if  $\gamma$  will be constructed by traversing through the nodes of the query tree using the pre-order method. The algorithm outputs  $\bar{c} = (\gamma(x_1), \dots, \gamma(x_k), b_1, \dots, b_\ell)$ . Note that  $\bar{c} < \bar{a}$  since by construction of  $\gamma$  follows that  $(\gamma(x_1), \dots, \gamma(x_k)) < (a_1, \dots, a_k)$ .

Let us consider the case  $\bar{a} \notin Q(D)$  and  $(D, \alpha) \models \varphi$ . If  $(a_{k+1}, \dots, a_{k+\ell}) > (b_1, \dots, b_\ell)$ , output the tuple  $\bar{c} = (a_1, \dots, a_k, b_1, \dots, b_\ell)$ . It is straightforward to see that  $\bar{c}$  is the sought tuple. If  $(a_{k+1}, \dots, a_{k+\ell}) < (b_1, \dots, b_\ell)$ , use the algorithm for the case that  $\bar{a} \notin Q(D)$  and  $(D, \alpha) \not\models \varphi$ , but the number  $\hat{p}$  is the number  $k$  (instead of setting  $\hat{p}$  to the minimum index such that  $[\alpha|_{\text{vpath}[x_p]}]$  is not in the  $x_p$ -list of  $[\alpha|_{\text{vpath}[x_p]}]$ ).

We now show that the choice of the tuple above is correct and that the tuple can be computed in time  $O(\text{poly}(Q)\|D\|)$ .

**Correctness** It is easy to see that the algorithm is correct if  $\bar{a} \in Q(D)$  or  $\bar{a}$  is smaller than the minimum tuple in  $Q(D)$ . The case that  $\bar{a} \notin Q(D)$  and  $(D, \alpha) \models \varphi$  and  $(a_{k+1}, \dots, a_{k+\ell}) < (b_1, \dots, b_\ell)$  is straightforward.

Let us consider the other cases: Let us first look upon the case that  $\bar{a} \notin Q(D)$  and  $(D, \alpha) \not\models \varphi$ . The constructed tuple  $\bar{c}$  belongs to the query result since the items  $[\gamma|_{\text{vpath}[v]}]$  are fit for all  $v \in \text{vars}(Q)$ . Let us assume for a contradiction that there is a tuple  $\bar{d} \in Q(D)$  with  $\bar{c} <_{\text{lex}} \bar{d} <_{\text{lex}} \bar{a}$  where  $\bar{c} := (\gamma(x_1), \dots, \gamma(x_p), b_1, \dots, b_\ell)$ . Then, there is a  $p' \in \{p, \dots, k + \ell\}$  such that for all  $j < p'$  is  $c_j = d_j$  and  $c_{p'} < d_{p'} \leq a_{p'}$ . (Here  $p$  is defined as in the algorithm above) Note that by construction is  $c_j = a_j$  for all  $j < p$  and, in particular, we have  $c_j = d_j = a_j$  for all  $j < p$ .

- *Case 1:  $p' = p$  and  $d_{p'} < a_{p'}$*

In that case, there is an item  $[\alpha|_{\text{vpath}[x_p]}] \cup \{x_p \mapsto d_{p'}\}$  in the  $x_p$ -list of  $[\alpha|_{\text{vpath}[x_p]}]$  with  $c_p < d_p < a_p$ . This violates the choice of  $c_p$ .

- *Case 2:  $p' = p$  and  $d_p = a_p$*

Since  $\bar{d} <_{\text{lex}} \bar{a}$  there is a  $\tilde{p}$  such that for all  $j < \tilde{p}$  it holds that  $d_j = a_j$  and  $d_{\tilde{p}} < a_{\tilde{p}}$ . Because  $\bar{d} \in Q(D)$  it follows that  $[\alpha|_{\text{vpath}[x_{\tilde{p}}]}] \cup \{x_{\tilde{p}} \mapsto d_{\tilde{p}}\}$  is in the  $x_{\tilde{p}}$ -list of  $[\alpha|_{\text{vpath}[x_{\tilde{p}}]}]$ .

Because  $\bar{d} \in Q(D)$  and since  $[\alpha|_{\text{vpath}[x_{\tilde{p}}]}]$  is not in the  $x_{\tilde{p}}$ -list of  $[\alpha|_{\text{vpath}[x_{\tilde{p}}]}]$  it holds that  $\tilde{p} \leq \hat{p}$  as  $d_{\tilde{p}} \neq a_{\tilde{p}}$  (Note that otherwise, the item  $[\alpha|_{\text{vpath}[x_{\tilde{p}}]}]$  must be fit)..

All in all, we have that there is an index  $\tilde{p} \in \{p+1, \dots, \hat{p}\}$  such that there is an item  $\iota$  in the  $x_{\tilde{p}}$ -list of  $[\alpha|_{\text{vpath}[x_{\tilde{p}}]}]$  with  $\alpha^\iota < a_{\tilde{p}}$ . This is a contradiction to the choice of  $p$ .

- *Case 3:  $p' > p$*

Then, there is an item  $[\gamma|_{\text{vpath}[x_{p'}]}] \cup \{x_{p'} \mapsto d_{p'}\}$  in the  $x_{p'}$ -list of  $[\gamma|_{\text{vpath}[x_{p'}]}]$  with  $c_{p'} < d_{p'}$ . But this is a contradiction to the fact that  $c_{p'}$  is the maximum element of the  $x_j$ -list of  $[\gamma|_{\text{vpath}[x_j]}]$ .

#### 4.8. Outputting the next smallest tuple

The proof for the case that  $\bar{a} \notin Q(D)$  and  $(D, \alpha) \models \varphi$  and  $(a_{k+1}, \dots, a_{k+\ell}) < (b_1, \dots, b_\ell)$  can be taken verbatim from the case that  $\bar{a} \notin Q(D)$  and  $(D, \alpha) \not\models \varphi$ , by removing the sentence “Because  $\bar{d} \in Q(D)$  and since  $[\alpha|_{\text{vpath}[x_{\hat{p}}]}]$  is not in the  $x_{\hat{p}}$ -list of  $[\alpha|_{\text{vpath}[x_{\hat{p}}]}]$  it holds that  $\tilde{p} \leq \hat{p}$  as  $d_p \neq a_p$  (Note that otherwise, the item  $[\alpha|_{\text{vpath}[x_{\hat{p}}]}]$  must be fit).”. Note that the fact  $\tilde{p} \leq \hat{p}$  is clear since  $\hat{p} = k$ .

**Running time** First, we compute the minimum element of  $Q(D)$  by starting the enumeration procedure until the first tuple will be output and test if  $\bar{a}$  is smaller than the minimum element. This can be done in  $\text{poly}(Q)$ . Then, we test in time  $\text{poly}(Q)$  if  $\bar{a} \in Q(D)$  and output in that case  $\bar{a}$ . In order to compute the number  $\hat{p}$ , it takes time  $\text{poly}(Q)$  to test if  $[\alpha|_{\text{vpath}[x_j]}]$  is in the  $x_j$ -list of  $[\alpha|_{\text{vpath}[x_j]}]$ , i.e., to test if  $[\alpha|_{\text{vpath}[x_j]}]$  is fit. Then, we find in time  $\text{poly}(Q)$  the maximum index  $p$  such that  $\tilde{p} \leq \hat{p}$  and  $[\alpha|_{\text{vpath}[x_{\tilde{p}}]}]$  is not the first item in the  $x_{\tilde{p}}$ -list of  $[\alpha|_{\text{vpath}[x_{\tilde{p}}]}]$ . To compute the element  $c_p$ , it takes time  $\text{poly}(Q) \log(\|D\|)$  to find the element using the binary search method if  $p = \hat{p}$ . Otherwise, one can simply take the previous item of the item  $\iota$  with  $a^t = \alpha(x_p)$ . The other components of  $\bar{c}$  can be computed in constant time. Thus, it takes time  $\text{poly}(Q) \log(\|D\|)$  to compute the tuple.

This concludes the proof of Lemma 4.31 (b).



## 5. Testing t-hierarchical conjunctive queries under updates

The aim of this section is to give a precise characterisation of the CQs for which *testing* can be done efficiently under updates. In this chapter, a proof for Theorem 3.6 is given.

Theorem 3.6 is restated in the following.

**Theorem 3.6** ([21]). *There is a dynamic algorithm that receives a  $t$ -hierarchical  $k$ -ary CQ  $Q$  and a  $\sigma$ -db  $D_0$ , and computes within  $t_p = \text{poly}(Q) \cdot O(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)$  and allows to test for an input tuple  $\bar{a} \in \mathbf{dom}^k$  if  $\bar{a} \in Q(D)$  within time  $t_t = \text{poly}(Q)$ , where  $D$  is the current database.*

The remainder of this chapter is based on [21].

Of course, according to Theorem 3.3 (a), the testing problem can be solved with constant update time and constant testing time for every  $q$ -hierarchical CQ. But the same holds true for the non- $q$ -hierarchical CQ  $Q_{S-E-T} := \{(x, y) : Sx \wedge Exy \wedge Ty\}$ . The corresponding dynamic algorithm simply uses 1-dimensional arrays  $\mathbf{A}_S$  and  $\mathbf{A}_T$  and a 2-dimensional array  $\mathbf{A}_E$  such that for all  $a, b \in \mathbf{dom}$  we have  $\mathbf{A}_E[a, b] = 1$  if  $(a, b) \in E^D$ , and  $\mathbf{A}_E[a, b] = 0$  otherwise, and  $\mathbf{A}_R[a] = 1$  if  $a \in R^D$ , and  $\mathbf{A}_R[a] = 0$  otherwise, for  $R \in \{S, T\}$ . When given an update command, the arrays can be updated within constant time. And when given a tuple  $(a, b) \in \mathbf{dom}^2$ , the **test** routine simply looks up the array entries  $\mathbf{A}_S[a]$ ,  $\mathbf{A}_E[a, b]$ ,  $\mathbf{A}_T[b]$  and returns the correct query result accordingly. To characterise the conjunctive queries for which testing can be done efficiently under updates, we consider  $t$ -hierarchical CQs (see Definition 3.4).

Obviously, it can be checked in time  $\text{poly}(Q)$  whether a given CQ  $Q$  is  $t$ -hierarchical. Note that every  $q$ -hierarchical CQ is  $t$ -hierarchical, and a *Boolean* query is  $t$ -hierarchical if and only if it is  $q$ -hierarchical. The aim of this chapter is to show Theorem 3.6, i.e., to show that the  $t$ -hierarchical CQs characterise CQs for which the *testing* problem can be solved efficiently under updates. In the remainder of this section, we give a proof for Theorem 3.6.

To avoid notational clutter, and without loss of generality, we restrict attention to queries  $Q_\varphi(u_1, \dots, u_k, b_1, \dots, b_\ell)$  where  $(u_1, \dots, u_k)$  is of the form  $(z_1, \dots, z_k)$  for pairwise distinct variables  $z_1, \dots, z_k$  and  $b_1, \dots, b_\ell \in \mathbf{dom}$ . We combine the array construction described above for the example query  $Q_{S-E-T}$  with the dynamic algorithm provided by Theorem 3.3 (a) and the following Lemma 5.1. To formulate the lemma, we need the following notation. A  $k$ -ary *generalised CQ* is of the form  $\{(z_1, \dots, z_k, b_1, \dots, b_\ell) : \varphi_1 \wedge \dots \wedge \varphi_m\}$  where  $k \geq 0$ ,  $z_1, \dots, z_k$  are pairwise distinct variables,  $m \geq 1$ ,  $\varphi_j$  is a conjunctive formula for each  $j \in [m]$ ,  $\text{free}(\varphi_1) \cup \dots \cup \text{free}(\varphi_m) =$

## 5. Testing $t$ -hierarchical conjunctive queries under updates

$\{z_1, \dots, z_k\}$ , and the quantified variables of  $\varphi_j$  and  $\varphi_{j'}$  are pairwise disjoint for all  $j, j' \in [m]$  with  $j \neq j'$  and disjoint from  $\{z_1, \dots, z_k\}$ . For each  $j \in [m]$  let  $\bar{z}^{(j)}$  be the sublist of  $\bar{z} := (z_1, \dots, z_k)$  that only contains the variables in  $\text{free}(\varphi_j)$  and let  $\bar{b}^{(j)}$  be the sublist of  $\bar{b} := (b_1, \dots, b_\ell)$  that only contains the constants in  $\varphi_j$ . I.e.,  $\bar{z}^{(j)}$  is obtained from  $\bar{z}$  by deleting all variables that do not belong to  $\text{free}(\varphi_j)$  and  $\bar{b}^{(j)}$  is obtained from  $\bar{b}$  by deleting all constants that do not appear in  $\varphi_j$ . Accordingly, for a tuple  $\bar{a} = (a_1, \dots, a_k) \in \mathbf{dom}^k$  by  $\bar{a}^{(j)}$  we denote the tuple that contains exactly those  $a_i$  where  $z_i$  belongs to  $\bar{z}^{(j)}$ . The query result of  $Q$  on a  $\sigma$ -db  $D$  is the set

$$Q(D) := \left\{ \bar{a}, \bar{b} \in \mathbf{dom}^k : D \models \varphi_j[\bar{a}^{(j)}] \text{ for each } j \in [m] \right\},$$

where  $D \models \varphi_j[\bar{a}^{(j)}]$  means that there is a homomorphism  $\beta_j : Q_j \rightarrow D$  for the query  $Q_j := \{ \bar{z}^{(j)}, \bar{b}^{(j)} : \varphi_j \}$ , with  $\beta_j(z_i) = a_i$  for every  $i$  with  $z_i \in \text{free}(\varphi_j)$  and  $\beta_j(c) = c$  for all  $c \in \text{cons}(\varphi_j)$ . For example,

$$Q'_{E-E-R} := \{ (x, y) : \exists v_1 Exv_1 \wedge \exists v_2 Eyv_2 \wedge \exists v_3 Rxyv_3 \}$$

is a generalised CQ that is equivalent to the CQ  $Q_{E-E-R}$ . The proof of the following lemma is given at the end of the section.

**Lemma 5.1.** *Every  $t$ -hierarchical CQ  $Q_\varphi(z_1, \dots, z_k, b_1, \dots, b_\ell)$  is equivalent to a generalised CQ  $Q' = \{ (z_1, \dots, z_k, b_1, \dots, b_\ell) : \varphi_1 \wedge \dots \wedge \varphi_m \}$  such that for each  $j \in [m]$  the CQ  $Q_j := \{ \bar{z}^{(j)} : \varphi_j \}$  is  $q$ -hierarchical or quantifier-free. Furthermore, there is an algorithm which decides in time  $\text{poly}(Q_\varphi)$  whether  $Q_\varphi$  is  $t$ -hierarchical, and if so, outputs an according  $Q'$ .*

When given a  $t$ -hierarchical CQ  $Q_\varphi(z_1, \dots, z_k)$ , use the algorithm provided by Lemma 5.1 to compute an equivalent generalised CQ  $Q'$  of the form  $\{ (z_1, \dots, z_k) : \varphi_1 \wedge \dots \wedge \varphi_m \}$  and let  $Q_j := \{ \bar{z}^{(j)}, \bar{b}^{(j)} : \varphi_j \}$  for each  $j \in [m]$ . W.l.o.g. assume that there is a  $m' \in \{0, \dots, m\}$  such that  $Q_j$  is  $q$ -hierarchical for each  $j \leq m'$  and  $Q_j$  is quantifier-free for each  $j > m'$ . We use in parallel, for each  $j \leq m'$ , the data structures provided by Theorem 3.3 (a) for the  $q$ -hierarchical CQ  $Q_j$ . In addition to this, we use an  $r$ -dimensional array  $\mathbf{A}_R$  for each relation symbol  $R \in \sigma$  of arity  $r := \text{ar}(R)$ , and we ensure that for all  $\bar{c} \in \mathbf{dom}^r$  we have  $\mathbf{A}_R[\bar{c}] = 1$  if  $\bar{c} \in R^D$ , and  $\mathbf{A}_R[\bar{c}] = 0$  otherwise. When receiving an update command  $\text{update } R(\bar{c})$ , we let  $\mathbf{A}_R[\bar{c}] := 1$  if  $\text{update} = \text{insert}$ , and  $\mathbf{A}_R[\bar{c}] := 0$  if  $\text{update} = \text{delete}$ , and in addition to this, we call the **update** routines of the data structure for  $Q_j$  for each  $j \leq m'$ . Upon input of a tuple  $(\bar{a}, \bar{d}) \in \mathbf{dom}^{k+\ell}$ , the **test** routine proceeds as follows. We test if  $\bar{d} = \bar{b}$ . If the test fails, we terminate the routine and output **False**. Otherwise we continue as follows. For each  $j \leq m'$ , it calls the **test** routine of the data structure for  $Q_j$  upon input  $\bar{a}^{(j)}$ . Additionally, it uses the arrays  $\mathbf{A}_R$  for all  $R \in \sigma$  to check if for each  $j > m'$  the quantifier-free query  $Q_j$  is satisfied by the tuple  $\bar{a}^{(j)}$ . All this is done within time  $\text{poly}(q)$ , and we know that  $\bar{a} \in Q(D)$  if and only if all these tests succeed. This completes the proof of Theorem 3.6.

For the remainder of this section we give a proof for Lemma 5.1.

*Proof of Lemma 5.1.* Along Definition 3.4 it is straightforward to construct an algorithm that decides in time  $\text{poly}(Q)$  whether a given CQ  $Q$  is t-hierarchical.

Let  $Q := Q_\varphi(z_1, \dots, z_k, b_1, \dots, b_\ell)$  be a given t-hierarchical CQ. Let  $A_0$  be the set of all atoms  $\psi$  of  $Q$  with  $\text{vars}(\psi) \subseteq \text{free}(Q)$ , and let  $\varphi_0$  be the quantifier-free conjunctive formula

$$\varphi_0 := \bigwedge_{\psi \in A_0} \psi.$$

For each  $Z \subseteq \text{free}(Q)$  let  $A_Z$  be the set of all atoms  $\psi$  of  $Q$  such that  $Z = \text{vars}(\psi) \cap \text{free}(Q)$  and  $\text{vars}(\psi) \not\supseteq Z$ . Let  $Z_1, \dots, Z_n$  (for  $n \geq 0$ ) be a list of all those  $Z \subseteq \text{free}(Q)$  with  $A_Z \neq \emptyset$ . For each  $j \in [n]$  let  $A_j := A_{Z_j}$  and let  $Y_j := (\bigcup_{\psi \in A_j} \text{vars}(\psi)) \setminus Z_j$ .

**Claim 5.2.**  $Y_j \cap Y_{j'} = \emptyset$  for all  $j, j' \in [n]$  with  $j \neq j'$ .

*Proof.* We know that  $Z_j \neq Z_{j'}$ . W.l.o.g. there is a  $z \in Z_j$  with  $z \notin Z_{j'}$ .

For contradiction, assume that  $Y_j \cap Y_{j'}$  contains some variable  $y$ . Then,  $y \in \text{vars}(\psi)$  for some  $\psi \in A_j$  and  $y \in \text{vars}(\psi')$  for some  $\psi' \in A_{j'}$ . By definition of  $A_j$  we know that  $\text{vars}(\psi) \cap \text{free}(Q) = Z_j$ , and hence  $z \in \text{vars}(\psi)$ . By definition of  $A_{j'}$  we know that  $\text{vars}(\psi') \cap \text{free}(Q) = Z_{j'}$ , and hence  $z \notin \text{vars}(\psi')$ . Hence,  $\psi \in \text{atoms}(z)$  and  $\psi' \notin \text{atoms}(z)$ . Since  $\psi \in \text{atoms}(y)$  and  $\psi' \in \text{atoms}(y)$ , we obtain that  $\text{atoms}(z) \cap \text{atoms}(y) \neq \emptyset$  and  $\text{atoms}(y) \not\subseteq \text{atoms}(z)$ . But by assumption,  $Q$  is t-hierarchical, and this contradicts condition (ii) of Definition 3.4.  $\square$

For each  $j \in [n]$  consider the conjunctive formula

$$\varphi_j := \exists y_1^{(j)} \dots \exists y_{\ell_j}^{(j)} \bigwedge_{\psi \in A_j} \psi,$$

where  $\ell_j := |Y_j|$  and  $(y_1^{(j)}, \dots, y_{\ell_j}^{(j)})$  is a list of all variables in  $Y_j$ . Using Claim 5.2, it is straightforward to see that

$$Q' := \{ (z_1, \dots, z_k, b_1, \dots, b_\ell) : \varphi_0 \wedge \bigwedge_{j \in [n]} \varphi_j \}$$

is a generalised CQ that is equivalent to  $Q$ . Furthermore,  $Q'$  can be constructed in time  $\text{poly}(Q)$ . To complete the proof of Lemma 5.1 we consider for each  $j \in [n]$  the CQ

$$Q_j := \{ \bar{z}^{(j)} : \varphi_j \},$$

where  $\bar{z}^{(j)}$  is a tuple of length  $|Z_j|$  consisting of all the variables in  $Z_j$ .

**Claim 5.3.**  $Q_j$  is q-hierarchical, for each  $j \in [n]$ .

*Proof.* First of all, note that  $Q_j$  satisfies condition (ii) of Definition 3.1, since  $\text{free}(Q_j) = Z_j$ ,  $\text{atoms}_{Q_j}(z) = A_j$  for every  $z \in Z_j$ , and  $\text{atoms}_{Q_j}(y) \subseteq A_j$  for every  $y \in Y_j = \text{vars}(Q_j) \setminus \text{free}(Q_j)$ .

For contradiction, assume that  $Q_j$  is not q-hierarchical. Then,  $Q_j$  violates condition (ii) of Definition 3.1. I.e., there are variables  $x_1, x_2 \in Z_j \cup Y_j$  and atoms

5. Testing  $t$ -hierarchical conjunctive queries under updates

$\psi_1, \psi_2, \psi_{12} \in A_j$  such that  $\text{vars}(\psi_1) \cap \{x_1, x_2\} = \{x_1\}$ ,  $\text{vars}(\psi_2) \cap \{x_1, x_2\} = \{x_2\}$ , and  $\text{vars}(\psi_{12}) \cap \{x_1, x_2\} = \{x_1, x_2\}$ . Since  $\text{vars}(\psi) \cap \text{free}(Q) = Z_j$  for all  $\psi \in A_j$ , we know that  $x_1, x_2 \notin \text{free}(Q)$ . Therefore,  $x, x' \in \text{vars}(q) \setminus \text{free}(Q)$ , and hence  $\psi_1, \psi_2, \psi_{12}$  are atoms of  $Q$  which witness that condition (i) of Definition 3.4 is violated. This contradicts the assumption that  $Q$  is  $t$ -hierarchical.  $\square$

This completes the proof of Lemma 5.1.  $\square$



## 6. Answering q-hierarchical conjunctive queries with aggregates under updates

This chapter is devoted to q-hierarchical conjunctive queries with aggregates. First of all, we provide an example that motivates queries with aggregates. In this chapter, Theorem 3.7, that is restated in the following, will be proven.

**Theorem 3.7.** *There is a dynamic algorithm that receives a q-hierarchical CQ with aggregates  $Q$  with aggregation time  $t_a$  and a  $\sigma$ -db  $D_0$ , and computes within  $t_i = O(1)$  initialisation time and  $t_p = \text{poly}(Q)t_a O(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)t_a$  and allows to*

- (a) *enumerate  $Q(D)$  with delay  $t_d = \text{poly}(Q)$ ,*
- (b) *test for an input tuple  $\bar{a} \in \text{dom}^k$  if  $\bar{a} \in Q(D)$  within time  $t_t = \text{poly}(Q)$ ,*
- (c) *enumerate the tuples that join and omit the result set  $Q(D)$  after an arbitrary number of updates on a previous chosen state the database received with delay  $t_{di} = \text{poly}(Q)$  and*
- (d) *compute the cardinality  $|Q(D)|$  in time  $t_c = O(1)$*

where  $D$  is the current database.

### 6.1. Examples for queries with aggregates

**Example 6.1.** *Let us assume we have a database  $D$  that consists of two tables  $\text{Salary}^D$  and  $\text{Person}^D$ . In  $\text{Salary}^D$  we have triples  $(a_{pid}, a_{project}, a_{salary})$  where  $a_{pid}$  is a personal id,  $a_{project}$  is the name of a project and  $a_{salary}$  is the salary that the person with id  $a_{pid}$  receives for project  $a_{project}$ . Here, we assume that a person might be involved in multiple projects. In  $\text{Person}^D$  we have pairs  $(a_{pid}, a_{name})$  where  $a_{pid}$  is a personal id and  $a_{name}$  is the name of the person with id  $a_{pid}$ . Let us consider the following query:*

$$Q = \left\{ (x_{pid}, x_{name}, x_{project}, x_{salary}) : \begin{array}{l} \text{Person}(x_{pid}, x_{name}) \wedge \\ \text{Salary}(x_{pid}, x_{project}, x_{salary}) \end{array} \right\}$$

*This query describes tuples  $(a_{pid}, a_{name}, a_{project}, a_{salary})$  where  $a_{pid}$  is the id for  $a_{name}$  who receives a salary of  $a_{salary}$  for project  $a_{project}$ . In the following, some interesting database queries are given.*

## 6. Answering $q$ -hierarchical conjunctive queries with aggregates under updates

- Output name and id of the people and the sum of their salaries.
- Output the average salary of each individual project.
- Output the project with highest total salary.

To treat such queries we now introduce queries with aggregates. Later, we will present queries for the statements in Example 6.1.

First, we provide a formal definition for aggregate functions that is similar to the definition for aggregate function in [73]. For every set  $X$  and every  $m \in \mathbb{N}_{\geq 0}$  and every tuple  $(x_1, \dots, x_m) \in X^m$  let  $\llbracket x_1, \dots, x_m \rrbracket$  be the multiset with the elements  $x_1, \dots, x_m$ . The number  $m$  is the cardinality, denoted by  $|\llbracket x_1, \dots, x_m \rrbracket|$ . With  $\llbracket X \rrbracket_m$  we denote the set of the multisets over  $X$  with cardinality  $m$ . For two multisets  $\llbracket x_1, \dots, x_m \rrbracket$  and  $\llbracket y_1, \dots, y_{m'} \rrbracket$  we define their union as

$$\llbracket x_1, \dots, x_m \rrbracket \cup \llbracket y_1, \dots, y_{m'} \rrbracket := \llbracket x_1, \dots, x_m, y_1, \dots, y_{m'} \rrbracket.$$

An *aggregate function*  $\mathcal{F}$  is a class of functions  $\mathcal{F} = \{f_n\}_{n \in \mathbb{N}}$  such that  $f_n : \llbracket X \rrbracket_n \rightarrow X$ .

**Example 6.2.** Examples for aggregate functions are:

- **sum** =  $\{s_n\}_{n \in \mathbb{N}}$  with  $X = \mathbb{N}$  and  $s_0 := 0$ , and  $s_n(\llbracket x_1, \dots, x_n \rrbracket) := x_1 + \dots + x_n$  for all  $n \in \mathbb{N}_{\geq 1}$ ,
- **prod** =  $\{p_n\}_{n \in \mathbb{N}}$  with  $X = \mathbb{N}$  and  $p_0 := 1$ , and  $p_n(\llbracket x_1, \dots, x_n \rrbracket) := x_1 \cdot \dots \cdot x_n$  for all  $n \in \mathbb{N}_{\geq 1}$ ,
- **count** =  $\{c_n\}_{n \in \mathbb{N}}$  with  $X = \mathbb{R}$  and  $c_0 := 0$ , and  $c_n(\llbracket x_1, \dots, x_n \rrbracket) := n$  for all  $n \in \mathbb{N}_{\geq 1}$ ,
- **avg** =  $\{a_n\}_{n \in \mathbb{N}}$  with  $X = \mathbb{N}$  and  $a_0 := 0$  and for all  $n \in \mathbb{N}_{\geq 1}$  is  $a_n(\llbracket x_1, \dots, x_n \rrbracket) := (x_1 + \dots + x_n)/n$ ,
- **max** =  $\{m_n\}_{n \in \mathbb{N}}$  with  $X = \mathbb{N} \cup \{-\infty\}$  and  $m_0 := -\infty$ , and for all  $n \in \mathbb{N}_{\geq 1}$  is  $m_n(\llbracket x_1, \dots, x_n \rrbracket) := \max\{x_1, \dots, x_n\}$ .

We often write  $fx$  instead of  $f(x)$  for all  $x \in \llbracket X \rrbracket$ .

In the remainder of this section we give examples for queries with aggregates. A formal definition of the syntax and semantics of queries with aggregates is given in Section 6.3.

To describe queries with aggregates, we introduce expressions that describes how to aggregate over the database elements. Let us assume, we have the query and the database from Example 4.4. Now, we want to define a query, that for every  $\bar{a} \in \varphi(D)$  we have the first three components of  $\bar{a}$  and a fourth component which is the sum over every corresponding fourth component in  $\varphi(D)$  for the first three components, i.e.

$$\left\{ (\alpha(y), \alpha(x_1), \alpha(x_2), n) : n = \sum_{\substack{\beta \supseteq \alpha|_{\{y, x_1, x_2\}}, \\ (D, \beta) \models \varphi}} \beta(x_3), (D, \alpha) \models \varphi \right\}. \quad (6.1)$$

### 6.1. Examples for queries with aggregates

$y$	$x_1$	$x_2$	$\text{sum}(x_3)$
1	1	4	1
1	1	5	2
1	1	6	7
1	2	4	1
1	2	5	2
1	2	6	7
1	3	4	1
1	3	5	2
1	3	6	7
2	4	2	13
2	8	2	13
2	9	2	13
3	2	1	1

$y$	$x_1$	$\text{max}(\text{prod}(x_2, \text{sum}(x_3)))$
1	1	42
1	2	42
1	3	42
2	4	26
2	8	26
2	9	26
3	2	1

$y$	$\text{count}(x_1)$	$\text{max}(\text{prod}(x_2, \text{sum}(x_3)))$
1	3	42
2	3	16
3	1	1

Table 6.1.: Enumeration of  $Q_{\text{sum}}(D_0)$ ,  $Q_{\text{max}}(D_0)$  and  $Q_{\text{count}}(D_0)$  on the database from Example 4.4.

To describe a query for the set (6.1), we can define

$$Q_{\text{sum}} := \{(y, x_1, x_2, \text{sum}(x_3)) : Eyx_1 \wedge Fyx_2x_3 \wedge Gyx_2x_3\}$$

The query result of  $Q_{\text{sum}}$  on the database of Example 4.4 is given in Table 6.1. Another example is a query that generates triple such that for every  $\bar{a} = (a_1, a_2, a_3, a_4) \in Q_{\text{sum}}(D)$ , we output the first component  $a_1$  and the second component  $a_2$  and as a third component the maximum of the product  $a_3$  and  $a_4$ , i.e.,

$$\left\{ (\alpha(y), \alpha(x_1), n) : n = \max_{\substack{\gamma \supseteq \alpha|_{\{y, x_1\}}, \\ (D, \gamma) \models \varphi}} \left\{ \gamma(x_2) \cdot \sum_{\substack{\beta \supseteq \gamma|_{\{y, x_1, x_2\}}, \\ (D, \beta) \models \varphi}} \beta(x_3) \right\}, (D, \alpha) \models \varphi \right\}. \quad (6.2)$$

A query for the set (6.2) can be defined as:

$$Q_{\text{max}} := \{(y, x_1, \text{max}(\text{prod}(x_2, \text{sum}(x_3)))) : Eyx_1 \wedge Fyx_2x_3 \wedge Gyx_2x_3\}.$$

Furthermore, we can create a query with `count` aggregates. For example, we can define a query  $Q_{\text{count}}$  that outputs, evaluated on a database  $D$ , all triples of the form  $(a_1, m, a_3)$  where the following condition holds.

- There is a  $a_2$  such that  $(a_1, a_2, a_3) \in Q_{\text{max}}$  and
- $m$  is the number of elements  $a_2 \in \text{adom}(D)$  with  $(a_1, a_2, a_3) \in Q_{\text{max}}$ ,

## 6. Answering $q$ -hierarchical conjunctive queries with aggregates under updates

i.e.,  $Q_{\text{count}}(D)$  is the set

$$\{(\alpha(x_1), m, \alpha(x_3)) : m = |\{\beta : \beta \supseteq \alpha|_{\{x_1, x_3\}}, (D, \beta) \models Q_{\text{max}}\}|, (D, \alpha) \models Q_{\text{max}}\}.$$

The query  $Q_{\text{count}}$  can be described as

$$Q_{\text{count}} := \{(y, \text{count}(x_1), \text{max}(\text{prod}(x_2, \text{sum}(x_3)))) : E y x_1 \wedge F y x_2 x_3 \wedge G y x_2 x_3\}.$$

See Table 6.1 for the query results of  $Q_{\text{sum}}$  and  $Q_{\text{max}}$  and  $Q_{\text{count}}$  on the database of Example 4.4. As we see in our examples all queries  $Q$  are of the form

$$\{(x_1, \dots, x_\ell, \text{expr}_1, \dots, \text{expr}_m, b_1, \dots, b_\ell) : \varphi\}$$

The enumeration results are grouped by  $x_1, \dots, x_p$ . The expressions  $\text{expr}_1, \dots, \text{expr}_m$  are expressions that contain aggregate functions, such as  $\text{sum}, \text{max}, \text{prod}, \text{count}$  and variables. We will later give a formal definition of syntax and semantics of *aggregate expressions*. The expressions will be defined such that for all expressions  $\text{expr}_j$  there exists a variable  $v_j \in \text{vars}(Q)$  such that the expression is inductively defined over the height of the variables in the subtree of  $T_{\text{free}}$  induced on  $\text{succ}(v_j)$  where  $T_{\text{free}}$  is the subtree of the  $q$ -tree of  $Q$  induced on  $\text{free}(Q)$ .

Note that the form of the queries depends on an existing  $q$ -tree of a query. A query that is not allowed is the following

$$Q_3 := \{(\text{count}(y), x_1, x_2, x_3) : E y x_1 \wedge F y x_2 x_3 \wedge G y x_2 x_3\}.$$

Such a query cannot be maintained under updates in our notion. Let us assume, there is a dynamic algorithm **A** that receives  $Q_3$  and a  $\sigma$ -db  $D_0$ , computes within arbitrary preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q) \cdot O(n^{1-\varepsilon})$  and allows to enumerate  $Q_3(D)$  with delay  $t_d = O(n^{1-\varepsilon})$  for all  $\varepsilon > 0$ . Such an algorithm enumerates the set

$$\{(n, \alpha_1(x_1), \alpha(x_2), \alpha(x_3)) : n = |\{\beta \supseteq \alpha|_{\{x_1, x_2, x_3\}} : (D, \beta) \models \varphi\}|, (D, \alpha) \models \varphi\}$$

where  $\varphi := E y x_1 \wedge F y x_2 x_3 \wedge G y x_2 x_3$ . If we delete the first component in every output tuple, we receive the set  $\{(\alpha_1(x_1), \alpha(x_2), \alpha(x_3)) : (D, \alpha) \models \varphi\}$ . This is the result set of the query  $Q' = \{(x_1, x_2, x_3) : \exists y \varphi\}$ . Note that we output these tuples without repetition. For a contradiction, let us assume that we output a tuple  $\bar{a} = (a_1, a_2, a_3)$  twice. Then, **A** enumerates  $(n, a_1, a_2, a_3)$  and  $(m, a_1, a_2, a_3)$  with  $n \neq m$ .

Note that by definition of  $Q_3(D)$ , both tuples can not belong to the set. All in all, we have an algorithm that computes within arbitrary preprocessing time a data structure that can be updated in sublinear time that allows to enumerate the set  $Q'(D)$  with sublinear delay. Since  $Q'$  is not  $q$ -hierarchical this is a contradiction to Theorem 3.13.

In the next section we give a formal model to evaluate aggregates under updates.

## 6.2. A model to evaluate aggregations

To evaluate aggregates under updates, we define the *aggregate implementation*. The model we introduce here relies on the concept of a MUD-algorithm described in [41].

## 6.2. A model to evaluate aggregations

An aggregate implementation is a triple  $(\xi, \circ, \eta)$  where  $\xi : \mathbf{dom} \rightarrow \mathcal{M}$  is a function that maps a database entry to a message,  $\circ : \mathcal{M} \times \mathcal{M} \rightarrow \mathcal{M}$  is a binary operation on  $\mathcal{M}$  such that  $(\mathcal{M}, \circ)$  is an Abelian group and  $\eta : \mathcal{M} \rightarrow \mathbf{dom}$  is a function that maps the message to a database entry. The main idea is that we use the aggregate implementation to evaluate aggregate functions in the following way. For a multiset  $\llbracket a_1, \dots, a_n \rrbracket$  we use the function  $\xi$  to map the database entries to messages, then we connect them with the binary function  $\circ$  and use the solution as input on  $\eta$  to receive an output. An aggregate implementation computes an aggregate function  $\{f_n\}_{n \in \mathbb{N}}$  if for all  $n \in \mathbb{N}_{\geq 1}$  and for all  $a_1, \dots, a_n \in \mathbf{dom}$  the following holds:

$$f_n(\llbracket a_1, \dots, a_n \rrbracket) = \eta(\xi(a_1) \circ \dots \circ \xi(a_n))$$

and  $f_0 = \eta(e)$  where  $e$  is the neutral element in the group  $(\mathcal{M}, \circ)$ .

**Example 6.3.** We consider now examples for aggregate implementations where  $\mathbf{dom} = \mathbb{N}$ :

- Let  $\mathcal{I}_{\text{sum}} := (\xi_{\text{sum}}, \circ_{\text{sum}}, \eta_{\text{sum}})$  with
  - $\mathcal{M} = \mathbb{Z}$ ,
  - $\xi_{\text{sum}}(a) := a$  for all  $a \in \mathbf{dom}$ ,
  - $\circ_{\text{sum}} = +$  and
  - $\eta_{\text{sum}}(a) = a$ .

$\mathcal{I}_{\text{sum}}$  computes **sum**. Note that  $\eta_{\text{sum}}(\xi_{\text{sum}}(a_1) \circ_{\text{sum}} \dots \circ_{\text{sum}} \xi_{\text{sum}}(a_n)) = a_1 + \dots + a_n \in \mathbb{N}$  if  $a_1, \dots, a_n \in \mathbb{N}$ .

- Let  $\mathcal{I}_{\text{prod}} := (\xi_{\text{prod}}, \circ_{\text{prod}}, \eta_{\text{prod}})$  with
  - $\mathcal{M} = \mathbb{Q}$ ,
  - $\xi_{\text{prod}}(a) := a$  for all  $a \in \mathbf{dom}$ ,
  - $\circ_{\text{prod}}$  is the multiplication operator of  $\mathbb{Q}$  and
  - $\eta_{\text{prod}}(a) = a$ .

$\mathcal{I}_{\text{prod}}$  computes **prod**. Note that  $\eta_{\text{prod}}(\xi_{\text{prod}}(a_1) \circ_{\text{prod}} \dots \circ_{\text{prod}} \xi_{\text{prod}}(a_n)) = a_1 \cdot \dots \cdot a_n \in \mathbb{N}$  if  $a_1, \dots, a_n \in \mathbb{N}$ .

- Let  $\mathcal{I}_{\text{count}} := (\xi_{\text{count}}, \circ_{\text{count}}, \eta_{\text{count}})$  with
  - $\mathcal{M} = \mathbb{Z}$ ,
  - $\xi_{\text{count}}(a) := 1$  for all  $a \in \mathbf{dom}$ ,
  - $\circ_{\text{count}} = +$  and
  - $\eta_{\text{count}}(a) = a$ .

$\mathcal{I}_{\text{count}}$  computes **count**.

- Let  $\mathcal{I}_{\text{avg}} := (\xi_{\text{avg}}, \circ_{\text{avg}}, \eta_{\text{avg}})$  with
  - $\mathcal{M} = \mathbb{Z}^2$ ,

## 6. Answering $q$ -hierarchical conjunctive queries with aggregates under updates

- $\xi_{\text{avg}}(a) := (a, 1)$  for all  $a \in \mathbf{dom}$ ,
- $(a_1, n_1) \circ_{\text{avg}} (a_2, n_2) = (a_1 + a_2, n_1 + n_2)$  for all  $(a_1, n_1), (a_2, n_2) \in \mathcal{M}$  and
- $\eta_{\text{avg}}((a, n)) = a/n$ .

$\mathcal{I}_{\text{avg}}$  computes **avg**.

- For the aggregate function **max** it is not sufficient to simply "remember" the maximum element of a list. If we remove every occurrence of the maximum element in the list during update steps, there is no information what the current maximum element is. To get rid of that problem, we use the following idea. To define an aggregate implementation  $(\xi_{\text{max}}, \circ_{\text{max}}, \eta_{\text{max}})$  for the aggregate function **max** we set  $\mathcal{M}$  as the set of functions that maps elements from  $\mathbb{N} \cup \{-\infty\}$  to  $\mathbb{Z}$ . The main idea is that for a multiset  $\llbracket a_1, \dots, a_n \rrbracket$  the current message is a function that maps every  $a \in \mathbb{N} \cup \{-\infty\}$  to the number of occurrences of  $a$  in  $\llbracket a_1, \dots, a_n \rrbracket$ . We call such a function the occurrence function of  $\llbracket a_1, \dots, a_n \rrbracket$ . For every  $a \in \mathbb{N} \cup \{-\infty\}$ , the value  $\xi_{\text{max}}(a)$  is the function  $f$  with  $f(a) = 1$  and  $f(b) = 0$  for all  $b \in (\mathbb{N} \cup \{-\infty\}) \setminus \{a\}$ . This is the occurrence function of  $\llbracket a \rrbracket$ . The binary operation  $\circ_{\text{max}}$  outputs on input of  $f, g \in \mathcal{M}$  the function  $u$  such that  $u(a) = f(a) + g(a)$  for all  $a \in \mathbb{N} \cup \{-\infty\}$ . In other words, if  $f$  is the occurrence function of  $\llbracket a_1, \dots, a_n \rrbracket$  and  $g$  the occurrence function of  $\llbracket b_1, \dots, b_m \rrbracket$ , then the function  $f \circ_{\text{max}} g$  is the occurrence function of  $\llbracket a_1, \dots, a_n, b_1, \dots, b_m \rrbracket$ .

It remains to show, that  $(\mathcal{M}, \circ_{\text{max}})$  is an Abelian group.

- $\circ_{\text{max}}$  is associative. Let  $p, q, r \in \mathcal{M}$ . Then, for all  $a \in \mathbb{N} \cup \{-\infty\}$  the following holds:  $[(p \circ_{\text{max}} q) \circ_{\text{max}} r](a) = [p \circ_{\text{max}} q](a) + r(a) = (p(a) + q(a)) + r(a) = p(a) + (q(a) + r(a)) = p(a) + [q \circ_{\text{max}} r](a) = [p \circ_{\text{max}} (q \circ_{\text{max}} r)](a)$  and, in particular, is  $(p \circ_{\text{max}} q) \circ_{\text{max}} r = p \circ_{\text{max}} (q \circ_{\text{max}} r)$ .
- The neutral element  $e \in \mathcal{M}$  is the function  $e(a) = 0$  for all  $a \in \mathbb{N} \cup \{-\infty\}$  since for all  $q \in \mathcal{M}$  and for all  $a \in \mathbb{N} \cup \{-\infty\}$  holds  $[q \circ_{\text{max}} e](a) = q(a) + e(a) = q(a)$ .
- For every  $q \in \mathcal{M}$  there is an inverse element  $q^{-1} \in \mathcal{M}$  where  $q^{-1}(a) = -q(a)$  for all  $a \in \mathbb{N} \cup \{-\infty\}$  since  $[q \circ_{\text{max}} q^{-1}](a) = q(a) + q^{-1}(a) = 0$  for all  $a \in \mathbb{N} \cup \{-\infty\}$ , i.e.,  $q \circ_{\text{max}} q^{-1} = e$ .
- $\circ_{\text{max}}$  is commutative. Let  $p, q \in \mathcal{M}$ . For all  $a \in \mathbb{N} \cup \{-\infty\}$  let  $[p \circ_{\text{max}} q](a) = p(a) + q(a) = q(a) + p(a) = [q \circ_{\text{max}} p](a)$ . In particular,  $p \circ_{\text{max}} q = q \circ_{\text{max}} p$ .

The post-processing function returns the maximum  $a \in \mathbb{N} \cup \{-\infty\}$  that appears in the multiset, i.e.,

$$\eta(q) := \begin{cases} \max \{a : a \in \mathbb{N} \cup \{-\infty\}, q(a) \neq 0\} & \text{if } q \neq e \\ -\infty & \text{otherwise} \end{cases}$$

It is straightforward to see that  $(\xi_{\text{max}}, \circ_{\text{max}}, \eta_{\text{max}})$  computes the aggregation function **max**.

## 6.2. A model to evaluate aggregations

For an aggregate implementation  $(\xi, \circ, \eta)$  a data structure that maintains the aggregation function under updates represents the current message  $m_{\text{current}} \in \mathcal{M}$  and the value  $\eta(m_{\text{current}})$ , i.e., the value  $\eta(m_{\text{current}}) = f_n(\llbracket a_1, \dots, a_n \rrbracket)$  if  $\llbracket a_1, \dots, a_n \rrbracket$  is the current multiset.

To initialise the data structure, we initialise  $m_{\text{current}}$  to the neutral element of  $(\mathcal{M}, \circ)$  and compute  $\eta(m_{\text{current}})$ . If we *insert* an element  $a$  to the multiset, we update  $m_{\text{current}}$  to the result of  $m_{\text{current}} \circ \xi(a)$  and compute the new value  $\eta(m_{\text{current}})$ . If we *delete* an element  $a$  from the multiset, we compute  $\xi(a)^{-1}$  (the inverse element of  $\xi(a)$  in the group  $(\mathcal{M}, \circ)$ ) and update  $m_{\text{current}}$  to the result of  $m_{\text{current}} \circ \xi(a)^{-1}$  and compute the new value  $\eta(m_{\text{current}})$ . It is straightforward to verify that  $\eta(m_{\text{current}})$  always represents the value  $f_n(\llbracket a_1, \dots, a_n \rrbracket)$ .

To design suitable data structures for  $\mathcal{I}_{\text{sum}}$  and  $\mathcal{I}_{\text{prod}}$  and  $\mathcal{I}_{\text{count}}$  and  $\mathcal{I}_{\text{avg}}$  it is sufficient to simply declare variables that store the values  $m_{\text{current}}$  and  $\eta(m_{\text{current}})$ .

For  $(\xi_{\text{max}}, \circ_{\text{max}}, \eta_{\text{max}})$  we represent  $m_{\text{current}}$  via an ordered list  $\mathcal{L}$  with the values  $\{(a, m_{\text{current}}(a)) : a \in \mathbb{N} \cup \{-\infty\}, m_{\text{current}}(a) \neq 0\}$ . The values are ordered descending by the first value of the tuples, i.e., if  $(a_2, b_2)$  is the successor element of  $(a_1, b_1)$  in  $\mathcal{L}$ , then  $a_1 > a_2$ . To lookup, insert and remove elements in the list fast, we use AVL-trees (see [65]). This takes time  $O(\log n)$  where  $n$  is the number of elements in the current multiset. To initialise the data structure, we simply initialise an empty list. This represents the neutral element of  $(\mathcal{M}, \circ_{\text{max}})$ . To update the data structure for  $m_{\text{current}}$  to  $m_{\text{current}} \circ_{\text{max}} p$  where  $p \in \xi(\mathbb{N} \cup \{-\infty\}) \cup \xi^{-1}(\mathbb{N} \cup \{-\infty\})$  let  $b \in \mathbb{N} \cup \{-\infty\}$  be the unique element with  $p(b) \in \{1, -1\}$ . Note that, since  $p \in \xi(\mathbb{N} \cup \{-\infty\}) \cup \xi^{-1}(\mathbb{N} \cup \{-\infty\})$ , it holds that  $p(a) = 0$  for all  $a \in (\mathbb{N} \cup \{-\infty\}) \setminus \{b\}$ . Then, we do the following steps. We lookup  $b$  in the AVL-tree. If it is not present, we insert  $(b, p(b))$  to  $\mathcal{L}$ . Otherwise, we modify  $(b, n)$  to  $(b, n + p(b))$ . If  $n + p(b) = 0$  remove  $(b, 0)$  from  $\mathcal{L}$ . The value of  $\eta(m_{\text{current}})$  is the first component of the first element in the list. It is straightforward to verify, that after an operation the list  $\mathcal{L}$  represents  $m_{\text{current}}$  and  $\eta(m_{\text{current}})$  is the maximum element  $a \in \mathbb{N} \cup \{-\infty\}$  with  $m_{\text{current}}(a) \neq 0$ . In particular, it takes time  $O(\log n)$  to compute  $m_{\text{current}} \circ_{\text{max}} p$  from  $m_{\text{current}}$  and for an element  $p \in \xi(\mathbb{N} \cup \{-\infty\}) \cup \xi^{-1}(\mathbb{N} \cup \{-\infty\})$ .

For an aggregate implementation  $(\xi, \circ, \eta)$  let  $t_\xi$  be the time it takes to compute  $\xi(a)$  for an  $a \in \mathbf{dom}$ ,  $t_{\xi^{-1}}$  be the time it takes to compute  $\xi^{-1}(a)$  for a  $a \in \mathbf{dom}$ ,  $t_o$  be the time to compute  $m_{\text{current}} \circ p$  if we receive  $m_{\text{current}}$  and  $p \in \xi(\mathbf{dom}) \cup \xi^{-1}(\mathbf{dom})$  as input and  $t_\eta$  the time it takes to compute  $\eta(m_{\text{current}})$ . For an aggregation function  $\{f_n\}_{n \in \mathbb{N}}$  let  $t_a(\{f_n\}_{n \in \mathbb{N}})$  be defined as follows. Let  $\mathcal{A}$  be the set of all aggregate implementations that compute  $\{f_n\}_{n \in \mathbb{N}}$ , then

$$t_a(\{f_n\}_{n \in \mathbb{N}}) := \min_{(\xi, \circ, \eta) \in \mathcal{A}} t_\xi + t_{\xi^{-1}} + t_o + t_\eta,$$

i.e., it is the time it takes to perform an update in the multiset for the aggregate.

For the aggregate functions in Example 6.2, it follows that  $t_{\xi_{\text{agg}}} = t_{\xi_{\text{agg}}^{-1}} = O(1)$  for  $\text{agg} \in \{\text{sum}, \text{prod}, \text{count}, \text{avg}, \text{max}\}$  and  $t_{\eta_{\text{agg}}} = t_{o_{\text{agg}}} = O(1)$  for  $\text{agg} \in \{\text{sum}, \text{prod}, \text{count}, \text{avg}\}$  and  $t_{\eta_{\text{max}}} = t_{o_{\text{max}}} \leq O(\log n)$ . In particular is  $t_a(\text{max}) = O(\log n)$  and  $t_a(\text{agg}) = O(1)$  for  $\text{agg} \in \{\text{sum}, \text{prod}, \text{count}, \text{avg}\}$ .

## 6. Answering $q$ -hierarchical conjunctive queries with aggregates under updates

We will now give syntax and semantics of aggregation expressions.

### 6.3. Syntax and semantics of aggregation expressions

In this section, we give a precise definition of the syntax and semantic of aggregation expressions. The notion relies on the notion of conjunctive queries with aggregates [31] but we will define aggregate expressions that relies on the  $q$ -tree of the  $q$ -hierarchical query. Let  $\text{AGGR}$  be the set of aggregate functions.

**Definition 6.4** (Syntax of aggregate expression). *Let  $Q$  be a  $q$ -hierarchical query and let  $T$  be a  $q$ -tree of  $Q$  and  $T_{\text{free}}$  be the induced subtree of  $T$  on  $\text{free}(Q)$ . For all  $v \in V$  we define the set  $\text{AGGEXPR}_v^{T_{\text{free}}}$  of aggregate expressions over  $v$ , inductively over the height of a node in  $T_{\text{free}}$ .*

- *Let  $v$  be a leaf in  $T_{\text{free}}$  and  $\mathcal{F} \in \text{AGGR}$ . Then,  $\mathcal{F}(v) \in \text{AGGEXPR}_v^{T_{\text{free}}}$ .*
- *Let  $v$  be a node of height  $\geq 1$  in  $T_{\text{free}}$  and  $\mathcal{F}, \mathcal{G} \in \text{AGGR}$  and let  $u_1, \dots, u_s$  be the children of  $v$  in  $T_{\text{free}}$  and  $\text{expr}_j \in \text{AGGEXPR}_{u_j}^{T_{\text{free}}}$  for all  $j \in [s]$ . Then,*
  - $\mathcal{G}(\mathcal{F}(\text{expr}_1, \dots, \text{expr}_s)) \in \text{AGGEXPR}_v^{T_{\text{free}}}$  and
  - $\mathcal{G}(\mathcal{F}(v, \text{expr}_1, \dots, \text{expr}_s)) \in \text{AGGEXPR}_v^{T_{\text{free}}}$  and
  - $\mathcal{G}(\langle v, \text{expr}_1, \dots, \text{expr}_s \rangle) \in \text{AGGEXPR}_v^{T_{\text{free}}}$

The last expression  $\mathcal{G}(\langle v, \text{expr}_1, \dots, \text{expr}_s \rangle)$  allows us to construct an aggregate that operates over the solution of all children of the current node (see Example 6.7).

If we consider the query  $\varphi$  from Example 4.4 and the corresponding  $q$ -tree  $T$ , and its induced subgraph  $T_{\text{free}}$  on  $\text{free}(Q)$  then  $\text{sum}(x_3) \in \text{AGGEXPR}_{x_3}^{T_{\text{free}}}$  and  $\text{count}(x_1) \in \text{AGGEXPR}_{x_1}^{T_{\text{free}}}$  and  $\text{max}(\text{prod}(x_2, \text{sum}(x_3))) \in \text{AGGEXPR}_{x_2}^{T_{\text{free}}}$ .

**Definition 6.5** ( $q$ -hierarchical query with aggregates). *A  $q$ -hierarchical query with aggregates is a conjunctive query of the form*

$$Q := \{(x_1, \dots, x_s, \text{expr}_1, \dots, \text{expr}_r, b_1, \dots, b_\ell) : \exists x_{k+1} \dots \exists x_m \varphi\}$$

*such that the query*

$$Q' = \{(x_1, \dots, x_k, b_1, \dots, b_\ell) : \exists x_{k+1} \dots \exists x_m \varphi\}$$

*with  $s \leq k$  is  $q$ -hierarchical and  $x_1, \dots, x_k \in \text{vars}$  and  $b_1, \dots, b_\ell \in \text{dom}$  and there is a  $q$ -tree  $T$  of  $Q'$  and an induced subtree  $T_{\text{free}}$  of  $T$  on  $\text{free}(Q')$  such that for all  $j \in [r]$  is  $\text{expr}_j \in \text{AGGEXPR}_{v_j}^{T_{\text{free}}}$  with  $v_j \in \text{free}(\varphi)$  and*

$$\{x_1, \dots, x_s\} \cup \bigcup_{j=1}^r \text{succ}(v_j) = \text{free}(Q').$$

*The  $q$ -tree of  $Q$  is the  $q$ -tree of  $Q'$  and  $\text{free}(Q) := \text{free}(Q')$ .*



### 6.3. Syntax and semantics of aggregation expressions

The queries  $Q_{\text{sum}}$ ,  $Q_{\text{max}}$  and  $Q_{\text{count}}$ , we considered at the beginning of this chapter, are examples for q-hierarchical queries with aggregates.

Now, we give a precise definition of the semantics of an aggregate expression.

**Definition 6.6** (Semantics of an aggregate expression). *Let  $Q$  be a q-hierarchical query with aggregates, let  $T$  be a q-tree of  $Q$  and let  $T_{\text{free}}$  be the induced subgraph of  $T$  on  $\text{free}(Q)$ .*

*For all  $\sigma$ -databases  $D$  and all  $v \in V(T_{\text{free}})$  and all  $\alpha : \text{vpath}[v] \rightarrow \text{adom}(D)$  let  $\llbracket \cdot \rrbracket^{(D, \alpha)} : \text{AGGEXPR}_v^{T_{\text{free}}} \rightarrow \mathbf{dom}$  be a function that is defined inductively over the height of  $v$  as follows:*

- *If  $v$  is a leaf in  $T_{\text{free}}$  and  $\mathcal{F} = \{f_n\}_{n \in \mathbb{N}} \in \text{AGGR}$ . Then,*

$$\llbracket \mathcal{F}(v) \rrbracket^{(D, \alpha)} := f_n \left\{ \langle a^\iota : \iota \in \mathcal{L}_v^{[\alpha]} \rangle \right\}$$

*The variable of  $\mathcal{F}(v)$  is  $v$ .*

- *If  $v$  is a node of height  $h \geq 1$  in  $T_{\text{free}}$  and let  $u_1, \dots, u_s$  be the children of  $v$  and  $\mathcal{F} = \{f_n\}_{n \in \mathbb{N}} \in \text{AGGR}$  and  $\mathcal{G} = \{g_n\}_{n \in \mathbb{N}} \in \text{AGGR}$ . Then,*

$$\begin{aligned} \llbracket \mathcal{G}(\mathcal{F}(\text{expr}_1, \dots, \text{expr}_s)) \rrbracket^{(D, \alpha)} &:= \\ g_n \left\{ f_s \left\{ \llbracket \text{expr}_j \rrbracket^{(D, \alpha^\iota)} : j \in [s] \right\} : \iota \in \mathcal{L}_v^{[\alpha]} \right\} \end{aligned}$$

$$\begin{aligned} \llbracket \mathcal{G}(\mathcal{F}(v, \text{expr}_1, \dots, \text{expr}_s)) \rrbracket^{(D, \alpha)} &:= \\ g_n \left\{ f_{s+1} \left( \langle a^\iota \rangle \cup \left\{ \llbracket \text{expr}_j \rrbracket^{(D, \alpha^\iota)} : j \in [s] \right\} \right) : \iota \in \mathcal{L}_v^{[\alpha]} \right\} \end{aligned}$$

$$\begin{aligned} \llbracket \mathcal{G}(\langle v, \text{expr}_1, \dots, \text{expr}_s \rangle) \rrbracket^{(D, \alpha)} &:= \\ g_n \left\{ \langle a^\iota, \llbracket \text{expr}_1 \rrbracket^{(D, \alpha^\iota)}, \dots, \llbracket \text{expr}_s \rrbracket^{(D, \alpha^\iota)} \rangle : \iota \in \mathcal{L}_v^{[\alpha]} \right\} \end{aligned}$$

*The variable of these expressions above is  $v$ .*

For all databases  $D$  and all queries with aggregates

$$Q = \{(x_1, \dots, x_m, \text{expr}_1, \dots, \text{expr}_r, b_1, \dots, b_\ell) : \varphi\}$$

is

$$Q(D) := \left\{ \begin{array}{l} (\alpha(x_1), \dots, \alpha(x_m), n_1, \dots, n_r, b_1, \dots, b_\ell) : \\ n_j = \llbracket \text{expr}_j \rrbracket^{(D, \alpha|_{\text{vpath}[v_j]})} \text{ for all } j \in [r] \text{ and } (D, \alpha) \models \varphi \end{array} \right\}.$$

A *single result* query with aggregates is a query of the form

$$Q = \{(\text{expr}) : \varphi\}$$

where  $\text{expr}$  is of the form  $\mathcal{F}(\text{expr}_1, \dots, \text{expr}_s)$  with  $\mathcal{F} \in \text{AGGR}$  and  $\text{expr}_j \in \text{AGGEXPR}_{u_j}^{T_{\text{free}}}$  for all  $j \in [s]$  and  $u_1, \dots, u_s$  are the children of  $v_{\text{root}}$  in  $T_{\text{free}}$ . For all  $\sigma$ -databases  $D$  let  $Q(D) := f_s \left\{ \llbracket \text{expr}_j \rrbracket^{(D, \emptyset)} : j \in [s] \right\}$ .

## 6. Answering $q$ -hierarchical conjunctive queries with aggregates under updates

The aggregation time  $t_a$  of a query with aggregates  $Q$  is defined as  $t_a(Q) := \max_{\mathcal{F} \in \mathcal{A}} t_a(\mathcal{F})$  where  $\mathcal{A}$  is the set of the aggregation function that occur in the expressions  $\text{expr}_1, \dots, \text{expr}_r$ . For example, for the query  $Q_{\text{count}}$  is  $t_a(Q_{\text{count}}) = O(\log(n))$  and for the query  $Q_{\text{sum}}$  is  $t_a(Q_{\text{sum}}) = O(1)$ .

**Example 6.7.** Recall the scenario from Example 6.1. We will now give queries for the statements given in the example:

- Output name and id of the people and the sum of their salaries:

$$Q_1 = \{(x_{\text{pid}}, x_{\text{name}}, \text{sum}(x_{\text{salary}})) : \exists x_{\text{project}} \varphi\}$$

where

$$\varphi := \text{Person}(x_{\text{pid}}, x_{\text{name}}) \wedge \text{Salary}(x_{\text{pid}}, x_{\text{project}}, x_{\text{salary}})$$

- Output the average salary of each individual project.

$$Q_2 = \{(x_{\text{project}}, \widetilde{\text{avg}}(\langle x_{\text{salary}}, \text{count}(x_{\text{pid}}) \rangle)) : \text{Salary}(x_{\text{pid}}, x_{\text{project}}, x_{\text{salary}})\}$$

where  $\widetilde{\text{avg}} = (a_n)_{n \in \mathbb{N}}$  is the aggregate

$$a_n(\langle x_1, y_1 \rangle, \dots, \langle x_n, y_n \rangle) := \frac{\sum_{i \in [n]} x_i y_i}{\prod_{i \in [n]} y_i}.$$

- Output the project with highest total salary:

$$Q_3 = \{(\text{max}_2(\langle x_{\text{project}}, \text{sum}(x_{\text{salary}}) \rangle)) : \exists x_{\text{pid}} \text{Salary}(x_{\text{pid}}, x_{\text{project}}, x_{\text{salary}})\}$$

where  $\text{max}_2 := \{m_n\}_{n \in \mathbb{N}}$  is defined as

$$m_n(\langle x_1, y_1 \rangle, \dots, \langle x_n, y_n \rangle) = y_i$$

where  $i$  is an index such that  $x_i = \max_{j \in [n]} x_j$ .

An aggregate implementation  $(\xi_{\widetilde{\text{avg}}}, \circ_{\widetilde{\text{avg}}}, \eta_{\widetilde{\text{avg}}})$  for  $\widetilde{\text{avg}}$  is the following: Let  $\mathcal{M} := \mathbb{Z} \times \mathbb{Q}$ . For every  $\langle x, y \rangle$  is  $\xi_{\widetilde{\text{avg}}}(\langle x, y \rangle) = (x \cdot y, y)$  and for all  $(z_1, q_1), (z_2, q_2) \in \mathbb{Z} \times \mathbb{Q}$  is  $(z_1, q_1) \circ_{\widetilde{\text{avg}}} (z_2, q_2) = (z_1 + z_2, q_1 * q_2)$ . The postprocessing function is  $\eta_{\widetilde{\text{avg}}}((z, q)) = z/q$ . It is easy to verify that the aggregate implementation  $(\xi_{\widetilde{\text{avg}}}, \circ_{\widetilde{\text{avg}}}, \eta_{\widetilde{\text{avg}}})$  computes  $\widetilde{\text{avg}}$ .

An aggregate implementation for  $\text{max}_2$  can be constructed very similar to the aggregate implementation for  $\text{max}$ .

In Section 6.4 we show how to compute the  $q$ -tree of a  $q$ -hierarchical conjunctive query with aggregates is  $q$ -hierarchical. In the Section 6.5 we give an application for queries with aggregates. We show how to modify a  $q$ -hierarchical conjunctive query  $Q$  to a Boolean  $q$ -hierarchical query  $Q_c$  with aggregates such that  $Q_c(D)$  is the number of tuples in the result set  $Q(D)$ , i.e.,  $Q_c(D) = |Q(D)|$ . Afterwards, we show, in the next section, how to enrich the data structure described in Section 4.1 to maintain queries with aggregates. Finally, we show how to reduce queries with aggregates to queries without aggregates such that we obtain routines like **enumerate**, **test** and **count** for queries with aggregates.

## 6.4. How to compute the q-tree of a q-hierarchical conjunctive query with aggregates

In this section we show how to compute the q-tree of a q-hierarchical conjunctive query with aggregates.

**Lemma 6.8.** *Let*

$$Q = \{(x_1, \dots, x_s, \text{expr}_1, \dots, \text{expr}_r, b_1, \dots, b_\ell) : \varphi\}$$

*be a query where  $\varphi$  is a conjunctive formula,  $\text{expr}_j$  for all  $j \in [r]$  is an aggregate expression,  $x_1, \dots, x_s \in \text{vars}$  and  $b_1, \dots, b_\ell \in \text{dom}$ . There is an algorithm that tests in time  $\text{poly}(Q)$  whether  $Q$  is a q-hierarchical CQ with aggregates and, if it is the case, computes the q-tree of  $Q$ .*

The remainder of this section is devoted to the proof of 6.8

Note that we defined the expressions over the q-tree of the aggregate expression free version  $Q'$ . Thus, we have to find out if there is a q-tree for the query

$$Q' = \{(x_1, \dots, x_s, \text{expr}_1, \dots, \text{expr}_r, b_1, \dots, b_\ell) : \varphi\}$$

that coincides with the expressions given in the query  $Q$ . In order to compute a corresponding q-tree we compute a tree  $T_{\text{expr}_q}$  for the aggregates expression  $\text{expr}_q$  for all  $q \in [r]$ . The tree  $T_{\text{expr}_q}$  is defined inductively.

- Let  $\text{expr}$  be an expression of the form  $\mathcal{F}(v)$ , then the tree  $T_{\text{expr}}$  is the tree that consists of the node  $v$  and has no edge, i.e.,  $V(T_{\text{expr}}) = \{v\}$  and  $E(T_{\text{expr}}) = \emptyset$ .
- Let  $\text{expr}$  be an expression of the form  $\mathcal{G}(\mathcal{F}(\text{expr}_1, \dots, \text{expr}_s))$ ,  $\mathcal{G}(v, \mathcal{F}(\text{expr}_1, \dots, \text{expr}_s))$  or  $\mathcal{G}(\langle \text{expr}_1, \dots, \text{expr}_s \rangle)$ . Then, the tree  $T_{\text{expr}}$  is the following.

$$V(T_{\text{expr}}) := V(T_{\text{expr}_1}) \cup \dots \cup V(T_{\text{expr}_s}) \cup \{v\}$$

$$E(T_{\text{expr}}) := E(T_{\text{expr}_1}) \cup \dots \cup E(T_{\text{expr}_s}) \cup \left\{ (v, v_j) : \begin{array}{l} v_j \text{ is the variable of} \\ \text{expr}_j \text{ for all } j \in [r] \end{array} \right\}$$

Note that it follows from the definition of aggregate expression that  $T_{\text{expr}_j}$  has to be a subtree in the q-tree of  $Q$  for all  $j \in [r]$ . In the following we define  $v_j$  as the variable of  $\text{expr}_j$  for all  $j \in [r]$ . If there are  $j, j' \in [r]$  with  $v_j \in \text{succ}(v_{j'})$ , then  $T_{\text{expr}_j} = T_{\text{expr}_{j'}}|_{\text{succ}(v_j)}$ . Clearly, this can be tested in  $\text{poly}(Q)$ . Let us now consider the set  $\Phi_j$  of atoms in  $\varphi$  where the variable  $v_j$  appears. We will show in the following that for all  $u \in \text{vars}(Q) \setminus V(T_{\text{expr}_j})$  with  $u \in \text{vars}(\Phi_j)$  we have that  $u$  occurs in every atom in  $\Phi_j$ . Assume for a contradiction that there is a  $u \in \text{vars}(Q) \setminus V(T_{\text{expr}_j})$  with  $\text{atoms}(v) \not\subseteq \text{atoms}(u)$ . Because of the fact that  $Q$  is q-hierarchical we have that  $\text{atoms}(u) \subset \text{atoms}(v)$  or  $\text{atoms}(u) \subseteq \text{atoms}(v) = \emptyset$ . Note that since  $u \in \text{vars}(\Phi_j)$  and  $v$  appears in every atom in  $\Phi_j$  we have that  $\text{atoms}(u) \subseteq \text{atoms}(v) = \emptyset$  is not possible. If

## 6. Answering $q$ -hierarchical conjunctive queries with aggregates under updates

$\text{atoms}(u) \subset \text{atoms}(v)$  it follows that  $u \in \text{succ}(v)$  and then  $u \in V(T_{\text{expr}_j})$  which violates  $u \in \text{vars}(Q) \setminus V(T_{\text{expr}_j})$ .

Let  $\tilde{Q}$  be the query

$$\{(w_1, \dots, w_p, b_1, \dots, b_\ell) : \tilde{\varphi}\}$$

we obtain if we modify  $Q$  in the following way.

- Remove every aggregate expression.
- Remove in  $\varphi$  every variable that belongs to  $\bigcup_{j=1}^r \text{succ}(v_j)$ . The conjunctive formula we obtain is  $\varphi'$ .
- $w_1, \dots, w_p$  are the remaining free variables in  $\tilde{\varphi}$ .

The algorithm for testing if  $Q$  is a valid  $q$ -hierarchical conjunctive query with aggregates works as follows.

1. For all  $j \in [r]$  compute the tree  $T_{\text{expr}_j}$ .
2. For all  $j, j' \in [r]$  test if  $T_{\text{expr}_j} = T_{\text{expr}_{j'}}|_{\text{succ}(v_j)}$  or  $V(T_{\text{expr}_j}) \cap T_{\text{expr}_{j'}} = \emptyset$ . If the result of one of these tests is false, return **false**.
3. Compute the query  $\tilde{Q}$  and test if  $\tilde{Q}$  is  $q$ -hierarchical. If the query is not  $q$ -hierarchical, return **false**.
4. Compute the  $q$ -tree  $\tilde{T}$  of  $\tilde{Q}$ .
5. For all  $j \in [r]$  such that there is no  $j' \in [r] \setminus \{j\}$  with  $v_{j'} \in \text{succ}(v_j) \setminus \{v_j\}$  do the following. If  $\text{vars}(\Phi_j) \setminus \text{succ}(v_j) = \emptyset$ , concatenate the tree  $T_{\text{expr}_j}$  to the root, i.e., add to  $\tilde{T}$  the nodes and the edges of  $T_{\text{expr}_j}$  and an edge from  $v_{\text{root}}$  to the root of  $T_{\text{expr}_j}$ . If  $\text{vars}(\Phi_j) \setminus \text{succ}(v_j) \neq \emptyset$ , choose the variable of  $u \in \text{vars}(\Phi_j) \setminus \text{succ}(v_j)$  with the lowest height in  $\tilde{T}$ . Add the tree  $T_{\text{expr}_j}$  as a child of  $v$ , i.e., we add to  $\tilde{T}$  the nodes and edges of  $T_{\text{expr}_j}$  and an edge from  $u$  to the root of  $T_{\text{expr}_j}$ .
6. Test if  $\tilde{T}$  is a  $q$ -tree for  $Q'$ . Then,  $Q$  is a valid  $q$ -hierarchical conjunctive query with aggregates and  $\tilde{T}$  the corresponding  $q$ -tree. Otherwise, the algorithm return **false**.

Note that every step in the algorithm can be done in  $O(\text{poly}(Q))$ . In the following claim we will show that the algorithm is correct.

**Claim 6.9.** *The algorithm returns true and a  $q$ -tree if and only if  $Q$  is a  $q$ -hierarchical conjunctive query with aggregates.*

*Proof.* The algorithm returns **true** if  $\tilde{T}$  (in the last step) is a  $q$ -tree for  $Q'$ . Note that  $T_{\text{expr}_j}$  is a subtree of  $\tilde{T}$ . Thus, the query  $Q$  is  $q$ -hierarchical and  $\tilde{T}$  is the corresponding  $q$ -tree.

Let us now consider the case that the algorithm returns **false**. If the algorithm returns **false** in Step 2, there are  $j, j' \in [r]$  with  $T_{\text{expr}_j} \neq T_{\text{expr}_{j'}}|_{\text{succ}(v_j)}$  and  $V(T_{\text{expr}_j}) \cap$

#### 6.4. How to compute the q-tree of a q-hierarchical conjunctive query with aggregates

$V(T_{\text{expr}_{j'}}) \neq \emptyset$ . It follows from the definition of aggregate expressions that  $T_{\text{expr}_j}$  and  $T_{\text{expr}_{j'}}$  can not appear in the same q-tree. Therefore, the query is not q-hierarchical.

If the algorithm returns **false** in Step 3 the query  $\tilde{Q}$  is not q-hierarchical. Then, there are two variables  $x, y \in \text{vars}(\tilde{Q})$  with  $\text{atoms}(x) \not\subseteq \text{atoms}(y)$ ,  $\text{atoms}(y) \not\subseteq \text{atoms}(x)$  and  $\text{atoms}(x) \cap \text{atoms}(y) \neq \emptyset$ . Then, we have that there is an atom  $\psi_x \in \text{atoms}(x) \setminus \text{atoms}(y)$  and an atom  $\psi_y \in \text{atoms}(y) \setminus \text{atoms}(x)$  and atom  $\psi_{xy} \in \text{atoms}(x) \cap \text{atoms}(y)$ . Query  $\tilde{Q}$  is obtained from  $Q$  by removing variables in the atom. Note that the variable that has been removed, are removed in every atom. Hence, there are atoms in  $Q$  with  $\psi_x \in \text{atoms}(x) \setminus \text{atoms}(y)$  and an atom  $\psi_y \in \text{atoms}(y) \setminus \text{atoms}(x)$  and atom  $\psi_{xy} \in \text{atoms}(x) \cap \text{atoms}(y)$ . Thus,  $Q'$  is not q-hierarchical.

If the algorithm returns **false** in step 6,  $\tilde{T}$  is not a valid q-tree. Assume for a contradiction that  $Q$  is q-hierarchical. Consider that case that  $\tilde{T}$  is not a q-tree since there is an atom  $\psi$  in  $Q$  where  $\text{vars}(\psi) \cup \{v_{\text{root}}\}$  does not form a path in  $\tilde{T}$ . Since  $Q$  is q-hierarchical there is a q-tree  $\hat{T}$  for  $Q'$ . In particular,  $\text{vars}(\psi) \cup \{v_{\text{root}}\}$  forms a path in  $\hat{T}$ . If there is a variable  $v \in \{v_1, \dots, v_r\} \cap \text{vars}(\psi)$  then by definition of the aggregate expressions, the path has the following form. Every node that appears before  $v$  in the path does not belong to  $V(T_{\text{expr}})$  (where  $\text{expr}$  is the expression where  $v$  is the variable of) and thus there is a path from  $v_{\text{root}}$  to  $u$  where  $u$  is the previous variable of  $v$  in the path. The path from  $v_{\text{root}}$  to  $u$  does also appear in the q-tree of  $\tilde{Q}$ . Then, it will be followed by  $v$ . Since  $T_{\text{expr}}$  must be appear identical in  $\tilde{T}$  and in  $\hat{T}$ , there must also be a path from  $v$  to the rest of the nodes in the path. This is a contradiction.

Let us consider the case that  $\tilde{T}$  violates (2) in Definition 4.1. Then  $\text{free}(Q) \neq \emptyset$  and there are at least two connected components of  $\tilde{T}_{\text{free}}$  where  $\tilde{T}_{\text{free}}$  is the induced subgraph of  $\tilde{T}$  of  $\text{free}(Q)$ . Without loss of generality, let  $T_1$  and  $T_2$  be two connected components and let  $T_1$  be the connected component with  $v_{\text{root}} \in V(T_1)$ .

Let  $u$  be the root of the subtree  $T_2$ . Let  $w$  be the parent node of  $u$  in  $\tilde{T}$ . If  $w \in \text{succ}(v_j)$  for a  $j \in [r]$  we have that a contradiction to the fact that  $Q$  is q-hierarchical since  $T_{\text{expr}_j}$  must be contained in the q-tree of  $Q$  and hence, the subtree of the q-tree of  $Q$  induced on  $\text{free}(Q)$  has two connected components. If  $w \notin \text{succ}(v_j)$  for a  $j \in [r]$ , i.e.,  $w \notin \bigcup_{j=1}^r \text{succ}(v_j)$  let us consider two cases.

- Case 1:  $u \notin \bigcup_{j=1}^r \text{succ}(v_j)$ . Since the q-tree of  $\tilde{Q}$  is a subgraph of  $\tilde{T}$  this is a contradiction to the fact that  $\tilde{Q}$  is q-hierarchical. Thus, **false** has to be output in step 3.
- Case 2:  $u \in \{v_1, \dots, v_r\}$ . Since  $u$  is free and  $\text{atoms}(u) \subseteq \text{atoms}(v)$  it follows that in every path from  $v_{\text{root}}$  to  $u$  there is the quantified variable  $v$ . Thus, there can no q-tree for  $Q$  exists and thus,  $Q$  is not q-hierarchical.

□

This concludes the proof of Lemma 6.8.

### 6.5. Counting the number of tuples in $Q(D)$

In this section we show how to use queries with aggregates to output the number of tuples in the result set  $Q(D)$ . To realise this, we define inductively over the height of the variables in the induced subtree  $T_{\text{free}}$  of a  $q$ -tree of  $Q$  on  $\text{free}(Q)$  the expression  $c_v$ .

**Definition 6.10** (Expression for Counting). *Let  $Q$  be a  $q$ -hierarchical query and let  $T_{\text{free}}$  be the induced subtree of a  $q$ -tree of  $Q$  on  $\text{free}(Q)$ . If  $v$  is a leaf in  $T_{\text{free}}$ , then  $c_v := \text{count}(v)$ . Otherwise, let  $u_1, \dots, u_\ell$  be the children of  $v$  in  $T_{\text{free}}$ . Then,  $c_v := \text{sum}(\text{prod}(c_{u_1}, \dots, c_{u_\ell}))$ .*

In the next lemma, we show that the expression  $c_v$  describes the number of tuples in the result.

**Lemma 6.11.** *Let*

$$Q = \{(x_1, \dots, x_k, b_1, \dots, b_\ell) : \varphi\}$$

*be a  $q$ -hierarchical query and let  $T_{\text{free}}$  be the subtree of a  $q$ -tree of  $Q$  induced on  $\text{free}(Q)$ . Let  $x_1, \dots, x_m$  be the children of  $v_{\text{root}}$  in  $T_{\text{free}}$ . The query*

$$Q_c := \{(\text{prod}(c_{x_1}, \dots, c_{x_m})) : \varphi\}$$

*describes the number of the tuples in  $\varphi(D)$ , i.e.,  $Q_c(D) = \{|Q(D)|\}$ .*

To prove Lemma 6.11, the following claims (Claim 6.12 and Claim 6.13) will be convenient.

**Claim 6.12.** *Let  $Q$  be a  $q$ -hierarchical query and let  $T_{\text{free}}$  be the induced subtree of the  $q$ -tree of  $Q$  on  $\text{free}(Q)$  and let  $D$  be a  $\sigma$ -db. Let  $\iota$  be an arbitrary fit item in the data structure for  $Q$  on  $D$ . Then the following equation holds.*

$$|\tilde{\mathcal{E}}^\iota| = \prod_{j=1}^{\ell} \left| \bigcup_{\tilde{\iota} \in \mathcal{L}_{w_j}^\iota} \tilde{\mathcal{E}}^{\tilde{\iota}} \right| \quad (6.3)$$

where  $w_1, \dots, w_\ell$  are the children of  $v^\iota$  in  $T_{\text{free}}$ .

*Proof.* From Lemma 4.12 we obtain:

$$\begin{aligned} |\tilde{\mathcal{E}}^\iota| &= \left| \left\{ \bigcup_{j=1}^{\ell} \alpha_j \cup \alpha^\iota : \begin{array}{l} \text{for all } w_j \in \text{child}(u) \cap \text{free}(Q) \text{ there is a } \tilde{\iota} \in \mathcal{L}_{w_j}^\iota \\ \text{such that } \alpha_j \in \tilde{\mathcal{E}}^{\tilde{\iota}} \end{array} \right\} \right| \\ &\stackrel{(+)}{=} \left| \left\{ (\alpha_1, \dots, \alpha_\ell) : \begin{array}{l} \text{for all } w_j \in \text{child}(u) \cap \text{free}(Q) \text{ there is a } \tilde{\iota} \in \mathcal{L}_{w_j}^\iota \\ \text{such that } \alpha_j \in \tilde{\mathcal{E}}^{\tilde{\iota}} \end{array} \right\} \right| \\ &\stackrel{(++)}{=} \left| \left\{ \alpha_1 \in \tilde{\mathcal{E}}^{\tilde{\iota}} : \tilde{\iota} \in \mathcal{L}_{w_1}^\iota \right\} \times \dots \times \left\{ \alpha_\ell \in \tilde{\mathcal{E}}^{\tilde{\iota}} : \tilde{\iota} \in \mathcal{L}_{w_\ell}^\iota \right\} \right| \\ &= \prod_{j=1}^{\ell} \left| \left\{ \alpha_j \in \tilde{\mathcal{E}}^{\tilde{\iota}} : \tilde{\iota} \in \mathcal{L}_{w_j}^\iota \right\} \right| = \prod_{j=1}^{\ell} \left| \bigcup_{\tilde{\iota} \in \mathcal{L}_{w_j}^\iota} \tilde{\mathcal{E}}^{\tilde{\iota}} \right|. \end{aligned}$$

### 6.5. Counting the number of tuples in $Q(D)$

In equation (+) we use the following fact. Every assignment  $\alpha \in \tilde{\mathcal{E}}^i$  can be identified with a tuple  $(\alpha_1, \dots, \alpha_\ell)$  where for all  $w_j \in \text{child}(u) \cap \text{free}(Q)$  there is a  $\tilde{\iota} \in \mathcal{L}_{w_j}^i$  and an  $\alpha_j \in \tilde{\mathcal{E}}^{\tilde{\iota}}$  with  $\alpha_j \subseteq \alpha$ , and, in particular, the number of assignments in  $\tilde{\mathcal{E}}^i$  is equal to the number of tuples  $(\alpha_1, \dots, \alpha_\ell)$  of the form described above. In (++) we rewrite the set to a Cartesian product.  $\square$

**Claim 6.13.** *Let  $Q$  be a  $q$ -hierarchical query and let  $T_{\text{free}}$  be the induced subtree of the  $q$ -tree of  $Q$  on  $\text{free}(Q)$  and let  $D$  be a  $\sigma$ -db. For all items  $i$  in the data structure of  $Q$  on  $D$  and for all  $u \in \text{child}(v^i) \cap \text{free}(Q)$  holds*

$$\llbracket c_u \rrbracket^{(D, \alpha^i)} = \left| \bigcup_{\iota \in \mathcal{L}_u^i} \tilde{\mathcal{E}}^\iota \right|.$$

*Proof.* We show this claim inductively over the height of an item in  $T_{\text{free}}$ .

For the induction base, let us consider an item  $i$  of height 1. Then, for all  $u \in \text{child}(v^i)$  the following holds.

$$\llbracket c_u \rrbracket^{(D, \alpha^i)} = \llbracket \text{count}(u) \rrbracket^{(D, \alpha^i)} = |\{a^\iota : \iota \in \mathcal{L}_u^i\}| \stackrel{(\star)}{=} \left| \bigcup_{\iota \in \mathcal{L}_u^i} \{a^\iota\} \right| \stackrel{(\star\star)}{=} \left| \bigcup_{\iota \in \mathcal{L}_u^i} \tilde{\mathcal{E}}^\iota \right|$$

The equation marked with  $(\star)$  follows from the fact that every element in the multiset appears once. The equations marked with  $(\star\star)$  follows from the fact that  $\tilde{\mathcal{E}}^i = \{a^\iota\}$  for all  $\iota \in \mathcal{L}_u^i$ .

For the inductive step, let us consider an item  $i$  of height  $> 1$ . For all  $u \in \text{child}(v^i)$  which are leafs in  $T_{\text{free}}$  the claim can be shown analogously to the induction base. Let us now consider a child  $u \in \text{child}(v^i)$  with height  $\geq 1$ . Let  $w_1, \dots, w_s$  be the children of  $u$  in  $T_{\text{free}}$ . Then, it follows that

$$\begin{aligned} \llbracket c_u \rrbracket^{(D, \alpha^i)} &= \llbracket \text{sum}(\text{prod}(c_{w_1}, \dots, c_{w_s})) \rrbracket^{(D, \alpha^i)} \\ &= s_n \left\{ p_s \left\{ \llbracket c_{w_j} \rrbracket^{(D, \alpha^\iota)} : j \in [s] \right\} : \iota \in \mathcal{L}_v^i \right\} \\ &= \sum_{\iota \in \mathcal{L}_v^i} \prod_{j=1}^s \llbracket c_{w_j} \rrbracket^{(D, \alpha^\iota)} \stackrel{(IH)}{=} \sum_{\iota \in \mathcal{L}_v^i} \prod_{j=1}^s \left| \bigcup_{\tilde{\iota} \in \mathcal{L}_{w_j}^\iota} \tilde{\mathcal{E}}^{\tilde{\iota}} \right| \\ &\stackrel{(\star)}{=} \sum_{\iota \in \mathcal{L}_v^i} |\tilde{\mathcal{E}}^\iota| \stackrel{(\star\star)}{=} \left| \bigcup_{\iota \in \mathcal{L}_u^i} \tilde{\mathcal{E}}^\iota \right| \end{aligned}$$

Equation  $(IH)$  follows from the induction hypothesis. The equation  $(\star\star)$  follows from the fact that  $\tilde{\mathcal{E}}^{\iota_1} \cap \tilde{\mathcal{E}}^{\iota_2} = \emptyset$  for two  $\iota_1, \iota_2 \in \mathcal{L}_u^i$  with  $\iota_1 \neq \iota_2$ , i.e., the union in the last line is a disjoint union. The equation  $(\star)$  follows from Claim 6.12.

This concludes the proof of Claim 6.13.  $\square$

## 6. Answering $q$ -hierarchical conjunctive queries with aggregates under updates

*Proof of Lemma 6.11.* Note that the number of tuples in  $\varphi(D)$  is equal to  $|\tilde{\mathcal{E}}^{[\emptyset]}|$ . The aim is to show that  $Q_c(D) = |\tilde{\mathcal{E}}^{[\emptyset]}|$ .

It follows that

$$\begin{aligned} Q_c(D) &= p_s \left\{ \llbracket c_{x_j} \rrbracket^{(D, \emptyset)} : j \in [s] \right\} = \prod_{j=1}^m \llbracket c_{x_j} \rrbracket^{(D, \emptyset)} \\ &\stackrel{(6.13)}{=} \prod_{j=1}^m \left| \bigcup_{\iota \in \mathcal{L}_{x_j}^i} \tilde{\mathcal{E}}^\iota \right| \stackrel{(6.12)}{=} |\tilde{\mathcal{E}}^{[\emptyset]}| \end{aligned}$$

This concludes the proof of Lemma 6.11.  $\square$

The next section shows how to enrich the data structure for  $q$ -hierarchical queries to maintain queries with aggregates.

### 6.6. The data structure for queries with aggregates

First of all, we give an example how to construct a data structure that maintains the query  $Q_{\text{count}}$  on the database given in Example 4.4. Recall the query  $Q_{\text{count}}$ :

$$Q_{\text{count}} := \{(y, \text{count}(x_1), \max(\text{prod}(x_2, \text{sum}(x_3)))) : E y x_1 \wedge F y x_2 x_3 \wedge G y x_2 x_3\}.$$

On every present item  $i$  in the data structure with  $v^i = x_2$ , we store two values  $w_{\text{list}}^{\text{sum}(x_3)}(i)$  and  $w_{\text{item}}^{\text{prod}(x_2, \text{sum}(x_3))}(i)$  where

$$w_{\text{list}}^{\text{sum}(x_3)}(i) := \sum_{\hat{i} \in \mathcal{L}_{x_3}^i} a^{\hat{i}} \quad \text{and} \quad w_{\text{item}}^{\text{prod}(x_2, \text{sum}(x_3))}(i) := a^i \cdot w_{\text{list}}^{\text{sum}(x_3)}(i),$$

i.e.,  $w_{\text{list}}^{\text{sum}(x_3)}(i)$  is the sum of the constants of the items in the  $x_3$ -list of  $i$  and  $w_{\text{item}}^{\text{prod}(x_2, \text{sum}(x_3))}(i)$  stores the product of the sum and its constant. Furthermore, we store for every item  $i$  with  $v^i = y$  a value  $w_{\text{list}}^{\text{count}(x_1)}(i)$  and  $w_{\text{list}}^{\max(\text{prod}(x_2, \text{sum}(x_3)))}(i)$ , where  $w_{\text{list}}^{\text{count}(x_1)}(i)$  is the number of items in the  $x_1$ -list of  $i$  and  $w_{\text{list}}^{\max(\text{prod}(x_2, \text{sum}(x_3)))}(i)$  is the maximum over the  $w_{\text{item}}^{\text{prod}(x_2, \text{sum}(x_3))}(\hat{i})$ -values for all items  $\hat{i}$  in the  $x_2$ -list of  $i$ , i.e.

$$w_{\text{list}}^{\text{count}(x_1)}(i) := |\mathcal{L}_{x_1}^i|$$

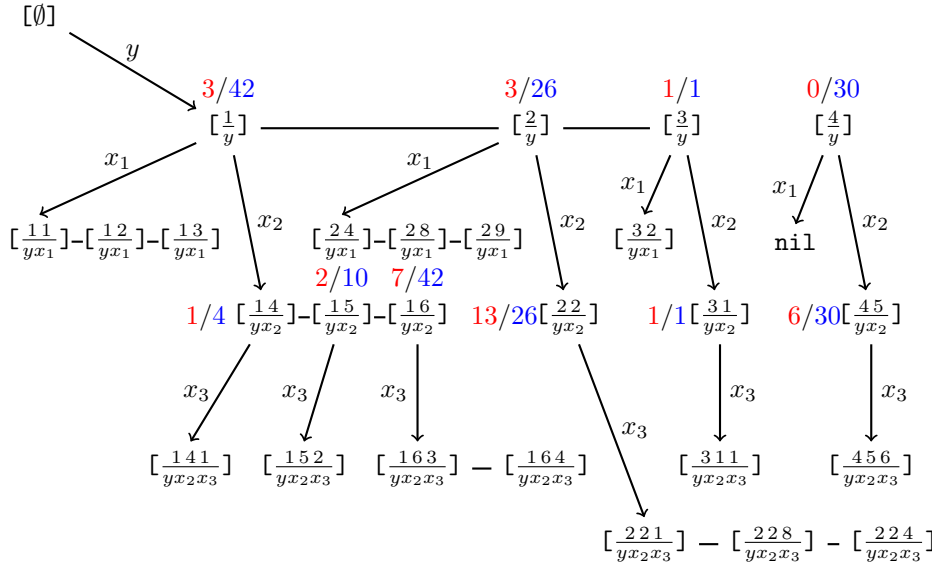
and

$$w_{\text{list}}^{\max(\text{prod}(x_2, \text{sum}(x_3)))}(i) := \max \left\{ w_{\text{item}}^{\text{prod}(x_2, \text{sum}(x_3))}(\hat{i}) : \hat{i} \in \mathcal{L}_{x_2}^i \right\}.$$

See Figure 6.1 for an illustration of the data structure together with the values. Every red (blue) number on an item  $i$  with  $v^i = x_2$  shows the value  $w_{\text{list}}^{\text{sum}(x_3)}(i)$  ( $w_{\text{item}}^{\text{prod}(x_2, \text{sum}(x_3))}(i)$ ), the red (blue) number above an item  $i$  with  $v^i = y$  shows the value  $w_{\text{list}}^{\text{count}(x_1)}(i)$  ( $w_{\text{list}}^{\max(\text{prod}(x_2, \text{sum}(x_3)))}(i)$ ).



Figure 6.1.: Illustration of Example 4.4 with aggregates.



Now, suppose that  $(4, 1)$  was inserted to  $E$ . As we know from our example the item  $[\frac{4}{y}]$  is going to be fit and we insert  $[\frac{41}{yx_1}]$  to the  $x_1$ -list of  $[\frac{4}{y}]$ . Therefore, we simply have to update the value  $w_{\text{list}}^{\text{count}(x_1)}(i)$  to 1. See Figure 6.2 for an illustration.

Different to the running example in Section 4.1 we now suppose that we delete  $(1, 6, 4)$  from  $F$ . As a consequence, the item  $[\frac{164}{y_{x_2x_3}}]$  loses his fit-status. Therefore, we subtract the value 4 from  $w_{\text{list}}^{\text{sum}(x_3)}([\frac{16}{y_{x_2}}])$  to obtain 3 and compute  $w_{\text{item}}^{\text{prod}(x_2, \text{sum}(x_3))}([\frac{1}{x_2}]) = 6 \cdot 3 = 18$ . Now, we have to update  $w_{\text{list}}^{\text{max}(\text{prod}(x_2, \text{sum}(x_3)))}([\frac{1}{y}])$ . We know that

$$w_{\text{list}}^{\max(\text{prod}(x_2, \text{sum}(x_3)))}(\lfloor \frac{1}{y} \rfloor) = \max \{4, 10, 42\}$$

To obtain the new value, we compute from  $\max \llbracket 4, 10, 42 \rrbracket$  the value  $\max \llbracket 4, 10, 42 \rrbracket \setminus \{42\} = \max \llbracket 4, 10 \rrbracket$  and then from  $\max \llbracket 4, 10 \rrbracket$  the value  $\max \llbracket 4, 10, 18 \rrbracket$ . Note that this can be done in  $O(\log n)$  since the aggregation time for  $\mathbf{max}$  is at most  $O(\log(n))$  where  $n := |\text{adom}(D)|$ .

See Figure 6.3 for an illustration after the update step.

Let us now go into technical details for arbitrary q-hierarchical queries with aggregates and databases. Let  $Q = \{(x_1, \dots, x_s, \text{expr}_1, \dots, \text{expr}_r) : \varphi\}$  be the input query. Let  $T_Q$  be a q-tree of  $Q$  and let  $T_{\text{free}}$  be the subtree of  $T_Q$  induced on  $\text{free}(Q)$ . Let  $D$  be a  $\sigma$ -db. For every  $v \in V(T')$  let  $\text{AE}(v)$  be the set of subexpressions that appear in one of the expressions  $\text{expr}_t$  for  $t \in [r]$  and  $\text{expr}_t$  is an expression with the variable  $v$ . To maintain queries with aggregates, we have to ensure that the following condition holds.

6. Answering  $q$ -hierarchical conjunctive queries with aggregates under updates

Figure 6.2.: Illustration of 4.4 with aggregates after  $(4, 1)$  was inserted to  $E$ .

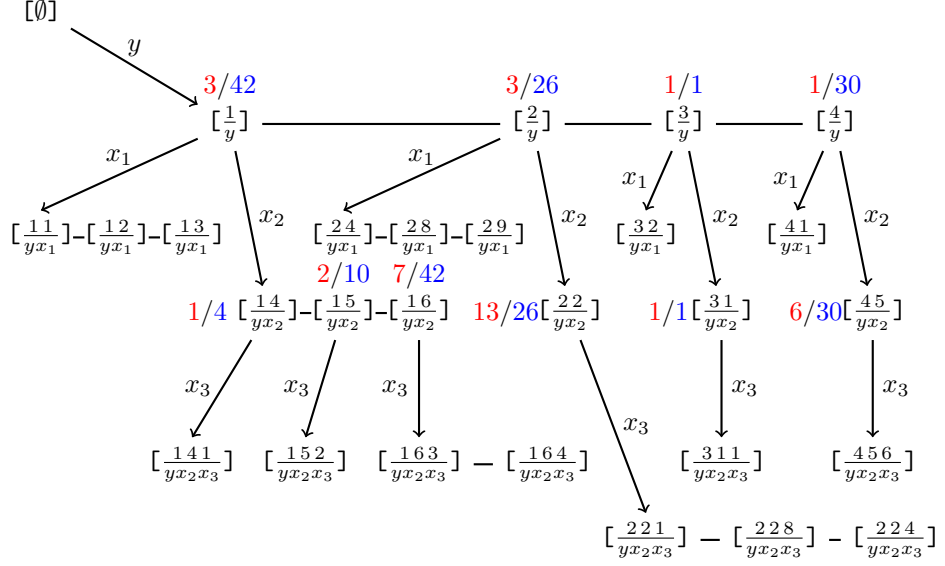
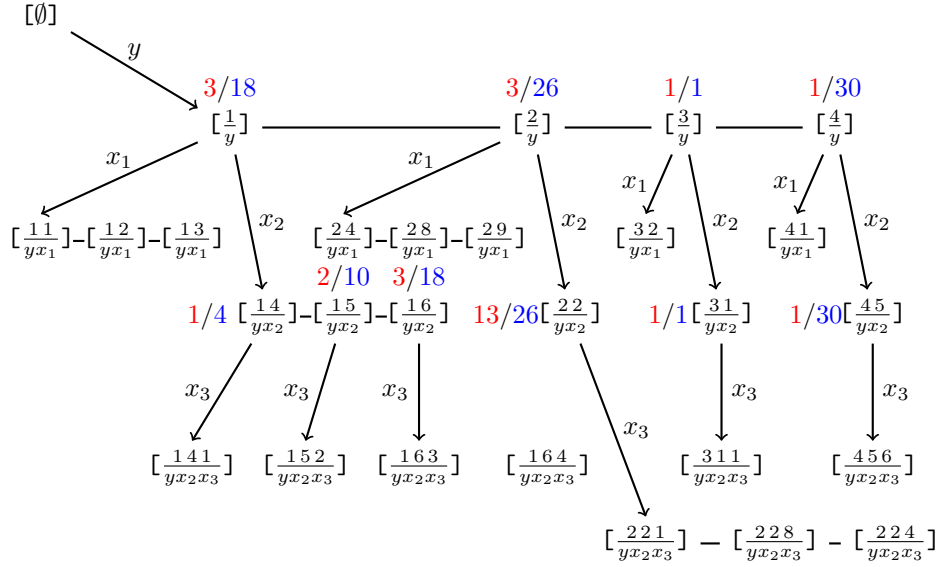


Figure 6.3.: Illustration of 4.4 with aggregates after  $(4, 1)$  was inserted to  $E$  and  $(1, 6, 4)$  deleted from  $F$ .



## 6.6. The data structure for queries with aggregates

**Condition 6.14.** For all present items  $i$  in the data structure for  $Q$  on  $D$ , where  $v^i$  is not a leaf in  $T_{\text{free}}$ , we store for every  $u \in \text{child}(v^i)$  and every  $\text{expr} \in \text{AE}(u)$  variables of the form  $w_{\text{list}}^{\text{expr}}(i)$  and  $w_{\text{item}}^{\text{expr}'}(i)$  (if  $\text{expr}$  is of the form  $\mathcal{G}(\text{expr}')$ ) such that the following is satisfied.

- if  $\text{expr}$  is of the form  $\mathcal{F}(u)$  then

$$w_{\text{list}}^{\mathcal{F}(u)}(i) = f_n \left\{ \{a^{\hat{i}} : \hat{i} \in \mathcal{L}_{u_j}^i\} \right\}$$

- and if  $\text{expr}$  is of the form  $\mathcal{G}(\text{expr}')$  then

$$w_{\text{list}}^{\mathcal{G}(\text{expr}')} (i) = g_n \left\{ \{w_{\text{item}}^{\text{expr}'}(\iota) : \iota \in \mathcal{L}_u^i\} \right\}$$

whereas if  $\text{expr}'$  is of the form  $\mathcal{F}(\text{expr}_1, \dots, \text{expr}_{\ell'})$  then

$$w_{\text{item}}^{\text{expr}'}(\iota) = f_{\ell'} \left\{ \{w_{\text{list}}^{\text{expr}_j}(\iota) : j \in [\ell']\} \right\},$$

and if  $\text{expr}'$  is of the form  $\mathcal{F}(u, \text{expr}_1, \dots, \text{expr}_{\ell'})$  then

$$w_{\text{item}}^{\text{expr}'}(\iota) = f_{\ell'+1} \left( \{a^{\iota}\} \cup \left\{ \{w_{\text{list}}^{\text{expr}_j}(\iota) : j \in [\ell']\} \right\} \right),$$

and if  $\text{expr}'$  is of the form  $\langle u, \text{expr}_1, \dots, \text{expr}_{\ell'} \rangle$  then

$$w_{\text{item}}^{\mathcal{G}(\text{expr}')}(\iota) = \langle a^{\iota}, w_{\text{list}}^{\text{expr}_1}(\iota), \dots, w_{\text{list}}^{\text{expr}_{\ell'}}(\iota) \rangle$$

The main idea is that these values are partial solutions of the aggregate values in the sense that  $w_{\text{list}}^{\text{expr}}(i) = \llbracket \text{expr} \rrbracket^{(D, \alpha^i)}$ . These values can be computed using a bottom-up algorithm. Such a value for  $i$  has to be changed only if there is an item in one of the  $u$ -lists of  $i$ , which value was changed, or the fit status of  $i$  changes. Therefore, whenever the fit status of an item changes, we can change all the partial solutions of  $w_{\text{list}}^{\text{expr}}(i)$  and, afterwards we recompute the partial solution for the parent item. We recompute the partial solutions of the corresponding parent item from the item we changed the values the last time, again and again, until we receive a parent item that has no aggregate or the item is the *start-item*.

In the following lemma, we show that the values stored in  $w_{\text{list}}^{\text{expr}}(i)$  are the partial solutions.

**Lemma 6.15.** Let  $Q$  be a  $q$ -hierarchical CQ with aggregates and  $D$  be a  $\sigma$ -db and  $T_{\text{free}}$  be the induced subgraph of the  $q$ -tree of  $Q$  on  $\text{free}(Q)$ . For every fit item  $i$  in the data structure for  $Q$  on  $D$  (that corresponds to  $T_Q$ ), of height  $\geq 1$  in  $T_{\text{free}}$  and every  $u \in \text{succ}(v^i)$  and every  $\text{expr} \in \text{AE}(u)$  the following holds:  $w_{\text{list}}^{\text{expr}}(i) = \llbracket \text{expr} \rrbracket^{(D, \alpha^i)}$ .

*Proof.* We prove this lemma by induction over the height of an item in  $T_{\text{free}}$ . For the induction base, let  $i$  be a fit item  $i$  of height 1 in  $T_{\text{free}}$ . Then, for every  $u \in \text{child}(v^i)$  has the expression  $\text{expr}$  the form  $\mathcal{F}(u)$ . Then,

$$w_{\text{list}}^{\mathcal{F}(u)}(i) = f_n \left\{ \{a^{\hat{i}} : \hat{i} \in \mathcal{L}_u^i\} \right\} = \llbracket \text{expr} \rrbracket^{(D, \alpha^i)}.$$

## 6. Answering $q$ -hierarchical conjunctive queries with aggregates under updates

For the inductive step, let us consider an item  $i$  of height  $h > 1$  and let  $u \in \text{child}(v^i)$  arbitrary and  $\text{expr}$  be an arbitrary expression in  $\text{AE}(u)$ . If  $\text{expr}$  is of the form  $\mathcal{F}(u)$ , the claim follows from the same argument as the induction base. If  $\text{expr}$  is of the form  $\mathcal{G}(\mathcal{F}(\text{expr}_1, \dots, \text{expr}_{s'}))$  where  $s' := |\text{child}(u)|$  then

$$\begin{aligned} w_{\text{list}}^{\text{expr}}(i) &= g_n \left\{ w_{\text{item}}^{\mathcal{F}(\text{expr}_1, \dots, \text{expr}_{s'})}(\iota) : \iota \in \mathcal{L}_u^i \right\} \\ &= g_n \left\{ f_{s'} \left\{ w_{\text{list}}^{\text{expr}_j}(\iota) : j \in [s'] \right\} : \iota \in \mathcal{L}_u^i \right\} \\ &\stackrel{(IH)}{=} g_n \left\{ f_{s'} \left\{ \llbracket \text{expr}_j \rrbracket^{(D, \alpha^i)} : j \in [s'] \right\} : \iota \in \mathcal{L}_u^i \right\} \\ &= \llbracket \text{expr} \rrbracket^{(D, \alpha^i)} \end{aligned}$$

Note that the equation marked with  $(IH)$  follows by the induction hypothesis. The case that  $\text{expr}$  is of the form  $\mathcal{G}(\mathcal{F}(u, \text{expr}_1, \dots, \text{expr}_{s'}))$  or  $\mathcal{G}(\langle u, \text{expr}_1, \dots, \text{expr}_{s'} \rangle)$  can be shown by an analogous way.  $\square$

To update the data structure we do the following. Every time we create a new item  $i$  during the update procedure, we set for all  $u \in \text{succ}(v^i)$  and for all  $\text{expr} \in \text{AE}(u)$  the value  $f_0$  if  $\text{expr}$  has the form  $\mathcal{F}(u)$  and the value  $g_0$  if  $\text{expr}$  has the form  $\mathcal{G}(\text{expr}')$ . Whenever the the fit status of an item  $i$  is changed we execute Algorithm 9.

---

**Algorithm 9** The fit status of  $\llbracket \frac{v_1, \dots, v_p}{b_1, \dots, b_p} \rrbracket$  has changed.

---

```

1: For all  $j \in \{0, \dots, p\}$  let  $\iota_j = \llbracket \frac{x_1, \dots, x_j}{b_1, \dots, b_j} \rrbracket$ .
2: for all  $\text{expr} \in \text{AE}(v_p)$  do
3:   if  $\text{expr}$  is of the form  $\mathcal{F}(v_p)$  then
4:      $w_{\text{list}}^{\text{expr}}(\iota_{p-1}) \leftarrow f_{n+1} \{a^{\hat{\iota}} : \hat{\iota} \in \mathcal{L}_{v_p}^{\iota_{p-1}}\}$ .
5:   if  $\text{expr}$  is of the form  $\mathcal{G}(\text{expr}')$  then
6:      $w_{\text{list}}^{\text{expr}}(\iota_{p-1}) = g_{n+1} \left\{ w_{\text{item}}^{\text{expr}'}(\hat{\iota}) : \hat{\iota} \in \mathcal{L}_{v_p}^{\iota_{p-1}} \right\}$ .
7: for all  $j = p-2$  to  $0$  do  $\triangleright$  Here  $v_0$  will be identified as  $v_{\text{root}}$ 
8:   for all  $\text{expr} \in \text{AE}(v_j)$  do  $\triangleright$   $\text{expr}$  is of the form  $\mathcal{G}(\text{expr}')$ 
9:     Let  $\text{expr}_q \in \text{AE}(v_{j+1})$  be the subexpression of  $\text{expr}$ .
10:    Let  $\ell' := |\text{child}(v_{j+1})|$ .
11:    if  $\text{expr}'$  is of the form  $\mathcal{F}(\text{expr}_1, \dots, \text{expr}_{\ell'})$  then
12:      Set  $w_{\text{item}}^{\text{expr}_q}(\iota_{j+1}) = f_{\ell'} \{w_{\text{list}}^{\text{expr}_s}(\iota_{j+1}) : s \in [\ell']\}$ .
13:    if  $\text{expr}'$  is of the form  $\mathcal{F}(v_{j+1}, \text{expr}_1, \dots, \text{expr}_{\ell'})$  then
14:      Set  $w_{\text{item}}^{\text{expr}_q}(\iota_{j+1}) = f_{\ell'} (\{a^{\iota_{j+1}}\} \cup \{w_{\text{list}}^{\text{expr}_s}(\iota_{j+1}) : s \in [\ell']\})$ .
15:    if  $\text{expr}'$  is of the form  $\langle v_{j+1}, \text{expr}_1, \dots, \text{expr}_{\ell'} \rangle$  then
16:      Set  $w_{\text{item}}^{\text{expr}_q}(\iota_{j+1}) = \langle a^{\iota_{j+1}}, w_{\text{list}}^{\text{expr}_1}(\iota_{j+1}), \dots, w_{\text{list}}^{\text{expr}_{\ell'}}(\iota_{j+1}) \rangle$ .
17:    Set  $w_{\text{list}}^{\text{expr}}(\iota_j) = g_n \left( \left\{ w_{\text{item}}^{\text{expr}_q}(\hat{\iota}) : \hat{\iota} \in \mathcal{L}_{v_{j+1}}^{\iota_{j+1}} \right\} \right)$ 

```

---

## 6.6. The data structure for queries with aggregates

**Lemma 6.16.** *Let  $Q$  be a  $q$ -hierarchical CQ with aggregates and let  $D$  be a  $\sigma$ -db. For all present items  $i$  in the data structure for  $Q$  on  $D$  the following holds.*

1. *After the data structure is initialised for the empty database, Condition 6.14 holds for  $i$  and*
2. *if Condition 6.14 holds for  $i$ , the condition still holds after the data structure was updated and Algorithm 9 was applied.*

*Proof.* When initialising the data structure for the empty database, we initialise a start item with empty lists and nothing else. In particular, Condition 6.14 holds.

Let us suppose that we have a data structure for a database  $D$  such that for all items Condition 6.14 holds for all present items in the data structure. Now the fit status of an item changed (or it is created or deleted). Then, it is straightforward to verify that the algorithm ensures for all  $j \in \{0, \dots, p\}$ , all possible  $w_{\text{list}}^{\text{expr}}(\iota_j)$  and all  $w_{\text{item}}^{\text{expr}}(\iota_j)$  are correct where  $\text{expr}$  are appropriate aggregate expressions.

Let us assume for a contradiction that there is an item  $\hat{i} \notin \{\iota_1, \dots, \iota_p\}$  where the  $w_{\text{list}}^{\text{expr}}(\hat{i})$  or the  $w_{\text{item}}^{\text{expr}}(\hat{i})$ -value is not correct and for all  $u \in \text{child}(v^i) \cap \text{free}(Q)$  and all  $\tilde{i} \in \mathcal{L}_u^i$  is the  $w_{\text{list}}^{\text{expr}}(\tilde{i})$  and the  $w_{\text{item}}^{\text{expr}}(\tilde{i})$  value is correct where  $\text{expr}$  are appropriate aggregate expressions. Then the corresponding multisets for  $w_{\text{list}}^{\text{expr}}(\hat{i})$  and  $w_{\text{item}}^{\text{expr}}(\hat{i})$  must be changed. But this can only happen if an item gets fit or unfit. Since  $\iota_p$  is the only item whose fit status was changed, this violates the fact that there is an item  $\hat{i} \notin \{\iota_1, \dots, \iota_p\}$  where the  $w_{\text{list}}^{\text{expr}}(\hat{i})$  or the  $w_{\text{item}}^{\text{expr}}(\hat{i})$ -value is not correct. In particular it follows that the values are correct for all present items. This concludes the proof.  $\square$

We now analyse the running time of Algorithm 9:

**Lemma 6.17.** *Let  $Q$  be a  $q$ -hierarchical CQ with aggregates and let  $t_a$  be the aggregation time of  $Q$ . Algorithm 9 takes time  $\text{poly}(Q)t_a$  on a data structure for an arbitrary  $\sigma$ -db  $D$ .*

*Proof.* Let us first consider line 4. The item  $\iota_p$  changes the fit status, therefore it is added to or removed from the  $v_p$ -list of  $\iota_{p-1}$ . Therefore, the value stored in  $w_{\text{list}}^{\text{expr}}(\iota_{p-1})$  (before the fit status of  $\iota_p$  changed) is equal to  $f_{n-1} \left\{ \left\{ a^{\hat{i}} : \hat{i} \in \mathcal{L}_{u_j}^i \setminus \{\iota_p\} \right\} \right\}$  or  $f_{n+1} \left\{ \left\{ a^{\hat{i}} : \hat{i} \in \mathcal{L}_{u_j}^i \cup \{\iota_p\} \right\} \right\}$ . To obtain the correct  $w_{\text{list}}^{\text{expr}}(\iota_{p-1})$  value update the value in time  $O(t_a(\mathcal{F}))$  by adding or removing a  $\iota_p$  to/from the multiset. With the same trick, we can update the value in line 6. Let us now consider line 12, 14 and 16). Since the values for the arguments  $a^{t_{j+1}}$  and  $w_{\text{item}}^{\text{expr}_s}(\tilde{i})$  for all  $s \in [\ell']$  are already stored and the number of these arguments are bounded by  $\ell' + 1$ , it takes at most  $(|\text{child}(v_{j+1})| + 1)t_a$  to compute  $w_{\text{item}}^{\text{expr}_q}(\iota_{j+1})$  from scratch. To update  $w_{\text{list}}^{\text{expr}}(\iota_j)$  in line 17 it takes at most  $2t_a$  to update the value, since we have to remove the "old"  $w_{\text{item}}^{\text{expr}_q}(\iota_{j+1})$  and insert the "new" value. Since the number of aggregate expressions is bounded by  $|Q|$  it follows for the running time of Algorithm 9:

$$|Q| \cdot 2t_a + \sum_{j=0}^{p-2} O(|Q|) [ (|\text{child}(v_{j+1})| + 1)t_a + 2t_a ] \leq \text{poly}(Q)t_a$$

□

All in all, we have shown that there is a data structure that maintains queries with aggregates under updates with update time  $t_u = \text{poly}(Q)t_a$ , since at most  $|\text{vars}(Q)|$  items change the fit status during an update.

## 6.7. A reduction to $q$ -hierarchical queries without aggregates

To obtain an enumeration algorithm for queries with aggregates, we define for every  $q$ -hierarchical CQ  $Q$  with aggregates a  $q$ -hierarchical CQ  $\tilde{Q}$  without aggregates and for every database  $D$  a database  $\tilde{D}$  such that  $Q(D) = \tilde{Q}(\tilde{D})$ . Then, we can use already known routines such as **enumerate**, **test** and **diff** for  $\tilde{Q}(\tilde{D})$  to obtain corresponding results for  $Q(D)$ . Furthermore, we guarantee that whenever the database  $D$  receives an update we have to update the data structure for  $\tilde{Q}(\tilde{D})$  very quickly (this will depend on the aggregation time of  $Q$ ).

Before we go into technical details, we show on an example how the reduction works. Let us recall the example from Section 6.6. We consider the query

$$Q_{\text{count}} := \{(y, \text{count}(x_1), \max(\text{prod}(x_2, \text{sum}(x_3)))) : E y x_1 \wedge F y x_2 x_3 \wedge G y x_2 x_3\}.$$

and the database from Example 4.4. See Figure 6.1 for an illustration of the data structure together with  $\text{count}(x_1)$  (the red values above the items  $i$  with  $v^i = y$ ) and  $\text{prod}(x_2, \text{sum}(x_3))$  (the blue values above the items  $i$  with  $v^i = y$ ).

Now from  $Q_{\text{count}}$ , we construct the following query by replacing the aggregate expressions  $\text{sum}(x_3)$  and  $\text{prod}(x_2, \text{sum}(x_3))$  by variables  $x_{\text{count}(x_1)}$  and  $x_{\text{prod}(x_2, \text{sum}(x_3))}$ . Furthermore we introduce for every aggregate expression relations  $R_{\text{prod}(x_2, \text{sum}(x_3))} y x_{\text{count}(x_1)}$  and  $R_{\text{count}(x_1)} y x_{\text{prod}(x_2, \text{sum}(x_3))}$ . We obtain the following  $q$ -hierarchical query:

$$\tilde{Q}_{\text{count}} := \left\{ (y, x_{\text{count}(x_1)}, x_{\text{prod}(x_2, \text{sum}(x_3))}) : \begin{array}{c} R_{\text{count}(x_1)} y x_{\text{count}(x_1)} \wedge \\ R_{\text{prod}(x_2, \text{sum}(x_3))} y x_{\text{prod}(x_2, \text{sum}(x_3))} \end{array} \right\}$$

of schema  $\tilde{\sigma} = \{R_{\text{count}(x_1)}, R_{\text{prod}(x_2, \text{sum}(x_3))}\}$ . We will now define a  $\tilde{\sigma}$ -db  $\tilde{D}$ . For every  $\text{expr} \in \{\text{count}(x_1), \text{prod}(x_2, \text{sum}(x_3))\}$  we define the binary relation  $R_{\text{expr}}$  as follows. There is a pair  $(a, b) \in R_{\text{expr}}(\tilde{D})$  if and only if the item  $i := [\frac{a}{y}]$  is present in the data structure for  $Q(D)$  (here, we mean the data structure described in Section 6.6) and  $\mathcal{L}_v^i \neq \emptyset$  where  $v$  is the variable of  $\text{expr}$  and  $b = w_{\text{list}}^{\text{expr}}(i)$ . In our example, we obtain the following database.

$$\begin{aligned} R_{\text{count}(x_1)}^{\tilde{D}} &= \{(1, 3), (2, 3), (3, 1)\} \\ R_{\text{prod}(x_2, \text{sum}(x_3))}^{\tilde{D}} &= \{(1, 42), (2, 26), (3, 1), (4, 30)\} \end{aligned}$$

See Figure 6.4 for an illustration of  $\tilde{Q}(\tilde{D})$ . A red dotted (blue dashed) arrow from an item  $i$  to an item  $\iota$  depicts that  $\iota$  is in the  $x_{\text{count}(x_1)}$ -list ( $x_{\text{prod}(x_2, \text{sum}(x_3))}$ -list) of  $i$ .

### 6.7. A reduction to $q$ -hierarchical queries without aggregates

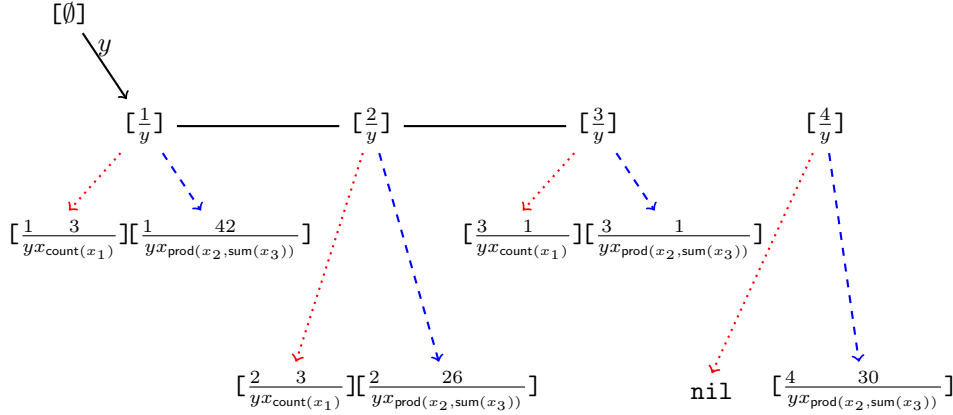


Figure 6.4.: Data structure for  $\tilde{Q}(\tilde{D})$ .

It is straightforward to see: When enumerating the tuples with the data structure in Figure 6.4, we get the result table shown in Table 6.1.

Let us assume that we insert  $(4, 1)$  to  $E$ . This forces that  $\mathcal{L}_{x_1}^{[\frac{4}{y}]}$  is not empty after the update in the data structure for  $\tilde{Q}$  on  $\tilde{D}$ . That is the reason why we have to insert the tuple  $(4, 1)$  to  $R_{\text{count}(x_1)}^{\tilde{D}}$ . We use the data structure from Section 6.6 to compute the new value  $w_{\text{list}}^{\text{count}(x_1)}([\frac{4}{y}])$  (see Figure 6.2 for the data structure after the update) and then we insert  $(4, 1)$  to  $R_{\text{count}(x_1)}^{\tilde{D}}$  using the known algorithm for updating data structures for  $q$ -hierarchical queries. See Figure 6.5 for an illustration of the data structure.

As a last update step in our example, we delete  $(1, 6, 4)$  from  $F$ . As a consequence, the value  $w_{\text{list}}^{\text{count}(x_1)}([\frac{1}{y}])$  changes from 42 to 18 (See Figure 6.3). Therefore, we delete the tuple  $(1, 42)$  from  $R_{\text{prod}(x_2, \text{sum}(x_3))}(\tilde{D})$  and insert the tuple  $(1, 18)$  to  $R_{\text{prod}(x_2, \text{sum}(x_3))}(\tilde{D})$ . See Figure 6.6 for the resulting data structure.

In the following lemma we will show that this idea works in general.

**Lemma 6.18.** *For every schema  $\sigma$ , every  $q$ -hierarchical query  $Q$  with aggregates of schema  $\sigma$  and every  $\sigma$ -db  $D$ , there is a schema  $\tilde{\sigma}$  and a  $q$ -hierarchical query  $\tilde{Q}$  without aggregates of schema  $\tilde{\sigma}$  and a  $\tilde{\sigma}$ -db  $\tilde{D}$  such that the following conditions hold:*

- (a)  $Q(D) = \tilde{Q}(\tilde{D})$  and
- (b)  $\tilde{Q}$  can be computed from  $Q$  in time  $O(\text{poly}(Q))$  and
- (c) an algorithm that maintains a data structure for  $\tilde{Q}$  on  $\tilde{\sigma}$ -dbs with initialisation time  $\tilde{t}_i$  and update time  $\tilde{t}_u$  can be used to obtain an algorithm for maintaining a data structure for  $Q$  on  $\sigma$ -dbs with initialisation time  $\tilde{t}_i$  and update time  $\text{poly}(Q)(\tilde{t}_u + t_a)$  where  $t_a$  is the aggregation time of  $Q$ .

6. Answering  $q$ -hierarchical conjunctive queries with aggregates under updates

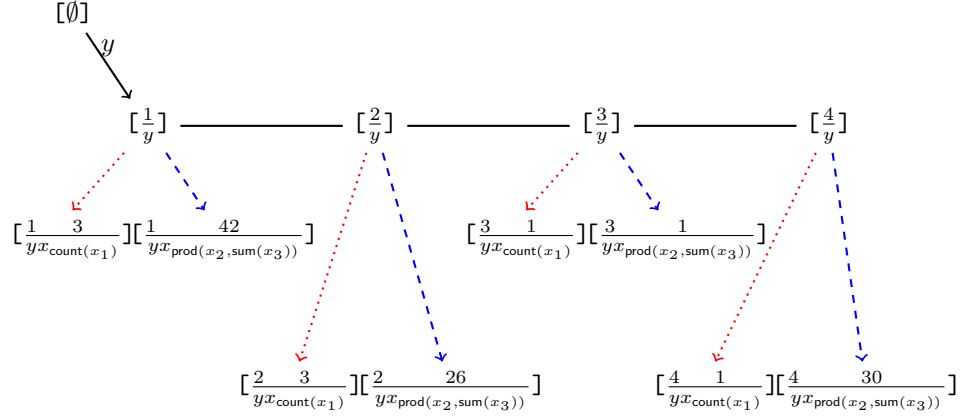


Figure 6.5.: Data structure for  $\tilde{Q}(\tilde{D})$  after inserting  $(4, 1)$  to  $E$ .

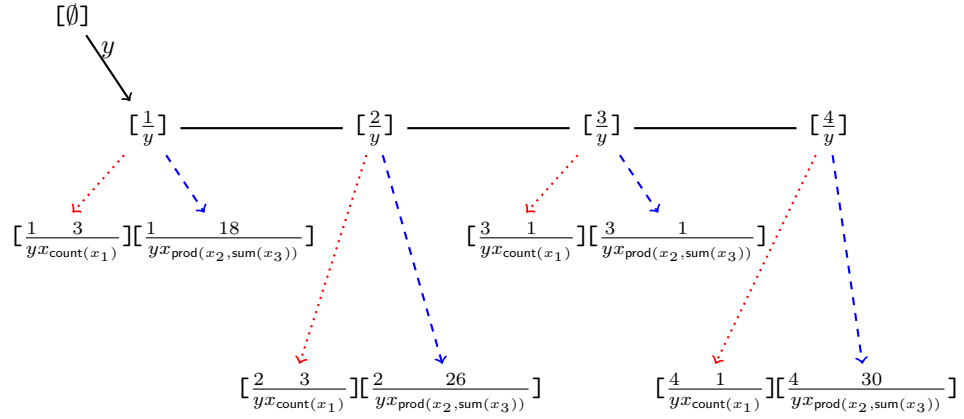


Figure 6.6.: Data structure for  $\tilde{Q}(\tilde{D})$  after inserting  $(4, 1)$  to  $E$  and deleting  $(1, 6, 4)$  from  $F$ .



## 6.7. A reduction to q-hierarchical queries without aggregates

*Proof.* Let

$$Q = \{(x_1, \dots, x_s, \text{expr}_1, \dots, \text{expr}_r, b_1, \dots, b_\ell) : \exists x_{k+1} \dots \exists x_m \varphi\}$$

be the input query and let

$$Q' = \{(x_1, \dots, x_k, b_1, \dots, b_\ell) : \exists x_{k+1} \dots \exists x_m \varphi\}$$

be the corresponding q-hierarchical query.

**Definition of  $\tilde{\sigma}$ ,  $\tilde{Q}$  and  $\tilde{D}$ .** The schema  $\tilde{\sigma}$  consists of relations  $R_p$  for every  $p \in [r]$  of arity  $\text{ar}(R_p) = |\text{vpath}[v_p]|$  where  $v_p$  is the variable of  $\text{expr}_p$  and of relations  $R \in \sigma$  such that there is an atom  $Ry_1 \dots y_{\text{ar}(R)}$  in  $Q$  such that  $\{y_1, \dots, y_{\text{ar}(R)}\} \cap \mathbf{var} \subseteq \{x_1, \dots, x_s, x_{k+1}, \dots, x_m\}$ . Let  $\tilde{\varphi}$  be the conjunction with the following atoms

- every atom  $\psi \in \text{atoms}(Q)$  where  $\text{vars}(\psi) \subseteq \{x_1, \dots, x_s, x_{k+1}, \dots, x_m\}$  is an atom in  $\tilde{\varphi}$  and
- for every  $p \in [r]$  there is an atom  $R_p(u_1, \dots, u_q, x_{\text{expr}_p})$  in  $\tilde{Q}$  where  $\{u_1, \dots, u_q\} = \text{vpath}[v_p]$  in which  $v_p$  is the variable of  $\text{expr}_p$ .

Let  $x_{j_1}, \dots, x_{j_q}$  be the variables in  $\{x_{k+1}, \dots, x_m\} \cap \text{vars}(\tilde{\varphi})$  and let  $c_1, \dots, c_p$  be the constants that appear in  $\text{atoms}(\varphi) \cap \text{atoms}(\tilde{\varphi})$ . The query  $\tilde{Q}$  is the following:

$$\tilde{Q} := \{(x_1, \dots, x_s, x_{\text{expr}_1}, \dots, x_{\text{expr}_r}, c_1, \dots, c_p) : \exists x_{j_1} \dots \exists x_{j_q} \tilde{\varphi}\}.$$

We now define the database  $\tilde{D}$ . For all  $R \in \sigma \cap \tilde{\sigma}$  let  $R(\tilde{D}) := R(D)$  and for all  $p \in [r]$  is

$$R_p(\tilde{D}) := \left\{ \begin{array}{l} (\alpha(u_1), \dots, \alpha(u_q), w_{\text{list}}^{\text{expr}}([\alpha])) : \\ \text{dom}(\alpha) = \{u_1, \dots, u_q\}, \text{ and} \\ [\alpha] \text{ is present in the data structure for } Q'(D), \text{ and} \\ \mathcal{L}_{v_p}^{[\alpha]} \neq \emptyset \text{ where } v_p \text{ is the variable of } \text{expr}_p \end{array} \right\}.$$

**The query  $\tilde{Q}$  is q-hierarchical.** The q-tree for  $\tilde{Q}$  can be constructed as follows. Let  $T$  be a q-tree for  $Q'$  and let  $T' = (V', E')$  be the subgraph of  $T$  induced on the set  $\{v_{\text{root}}, x_1, \dots, x_s, x_{j_1}, \dots, x_{j_q}\}$ . Note that by definition of queries with aggregates it follows that  $T'$  is connected. The tree  $\tilde{T} = (\tilde{V}, \tilde{E})$  is defined as follows

$$\tilde{V} := V' \cup \{x_{\text{expr}_p} : p \in [r]\}$$

$$\tilde{E} := E' \cup \{(w, x_{\text{expr}_p}) : w \text{ is the parent node of the variable of } x_{\text{expr}_p} \text{ in } T, p \in [r]\}$$

Clearly  $\tilde{V} = \text{vars}(\tilde{Q}) \cup \{v_{\text{root}}\}$ . It remains to show that for every atom  $\psi \in \text{atoms}(\tilde{Q})$  the set  $\{v_{\text{root}}\} \cup \text{vars}(\psi)$  forms a path in  $\tilde{T}$ . If  $\psi \in \text{atoms}(Q)$ , then  $\text{vars}(\psi) \subseteq \{x_1, \dots, x_s, x_{j_1}, \dots, x_{j_q}\}$  and  $\{v_{\text{root}}\} \cup \text{vars}(\psi)$  forms a path in  $T$ . Since  $T'$  is the induced subgraph of  $T$  on  $x_1, \dots, x_s, x_{j_1}, \dots, x_{j_q}$  and  $T'$  is a subgraph of  $\tilde{T}$ , it follows

## 6. Answering $q$ -hierarchical conjunctive queries with aggregates under updates

that  $\{v_{root}\} \cup \text{vars}(\psi)$  forms a path in  $\tilde{T}$ . If  $\psi \notin \text{atoms}(\tilde{Q})$ , there is a  $p \in [r]$  such that  $\psi$  is of the form  $R_p(u_1, \dots, u_q, x_{\text{expr}_p})$  where  $\{u_1, \dots, u_q\} = \text{vpath}[v_p]$  in which  $v_p$  is the variable of  $\text{expr}_p$ . Since  $\{u_1, \dots, u_q\} = \text{vpath}[v_p]$  the set  $\{v_{root}, u_1, \dots, u_q\}$  forms a path in  $T'$  and since  $(u_q, x_{\text{expr}_p}) \in \tilde{E}$ , it follows that  $\{v_{root}\} \cup \text{vars}(\psi)$  forms a path in  $\tilde{T}$ . Clearly, each of these paths starts from  $v_{root}$ . Since  $T'$  is a induced subgraph on  $x_1, \dots, x_s, x_{j_1}, \dots, x_{j_q}$  and it is a subgraph of  $T$ , the set  $x_1, \dots, x_s$  is connected. Since for every  $p \in [r]$  there is an edge  $(y, x_{\text{expr}_p})$  for a  $y \in \{x_1, \dots, x_s\}$ , the set  $x_1, \dots, x_s, x_{\text{expr}_1}, \dots, x_{\text{expr}_q}$  is connected in  $\tilde{T}$ . Therefore, the query  $\tilde{Q}$  is  $q$ -hierarchical.

**Proof of part (a).** As an abbreviation, we set  $\alpha$  be an assignment such that  $\bar{a} = (\alpha(x_1), \dots, \alpha(x_s))$  and  $\bar{c} := (c_1, \dots, c_p)$ .

$$\begin{aligned}
\bar{s} \in Q(D) &\Leftrightarrow \bar{s} = (\bar{a}, n_1, \dots, n_r, \bar{b}) \text{ where } n_j = \llbracket \text{expr}_j \rrbracket^{(D, \alpha|_{\text{vpath}[v_j]})} \\
&\text{for all } j \in [r], \text{ and there is a } \beta \supseteq \alpha \\
&\text{with } (D, \beta) \models \varphi \text{ and } \text{dom}(\beta) = \text{vars}(Q') \\
&\Leftrightarrow \bar{s} = (\bar{a}, w_{\text{list}}^{\text{expr}_1}(\llbracket \alpha|_{\text{vpath}[v_1]} \rrbracket), \dots, w_{\text{list}}^{\text{expr}_r}(\llbracket \alpha|_{\text{vpath}[v_r]} \rrbracket), \bar{b}) \\
&\text{and there is a } \beta \supseteq \alpha \text{ with} \\
&(D, \beta) \models \psi \text{ for all } \psi \in \text{atoms}(Q') \text{ and } \text{dom}(\beta) = \text{vars}(Q') \\
&\stackrel{(\star)}{\Leftrightarrow} \bar{s} = (\bar{a}, w_{\text{list}}^{\text{expr}_1}(\llbracket \alpha|_{\text{vpath}[v_1]} \rrbracket), \dots, w_{\text{list}}^{\text{expr}_r}(\llbracket \alpha|_{\text{vpath}[v_r]} \rrbracket), \bar{c}) \\
&\text{and } (D, \alpha) \models \psi \text{ for all } \psi \in \text{atoms}(Q) \cap \text{atoms}(\tilde{Q}) \\
&\text{and for all } R_p(u_1, \dots, u_t, x_{\text{expr}_p}) \in \text{atoms}(\tilde{Q}) \\
&\text{holds } (\alpha(u_1), \dots, \alpha(u_t), w_{\text{list}}^{\text{expr}_p}(\llbracket \alpha|_{\text{vpath}[v_r]} \rrbracket)) \in R_p(\tilde{D}) \\
&\Leftrightarrow \bar{s} \in \tilde{Q}(\tilde{D})
\end{aligned}$$

To prove that equation  $(\star)$  is correct, we show that there is a  $\beta \supseteq \alpha$  such that  $(D, \beta) \models \psi$  for all  $\psi \in \text{atoms}(Q')$  if and only if  $(D, \alpha) \models \psi$  for all  $\psi \in \text{atoms}(Q') \cap \text{atoms}(\tilde{Q})$  and for all atoms of the form  $R_p(u_1, \dots, u_r, x_{\text{expr}_p})$  we have  $(\alpha(u_1), \dots, \alpha(u_r), w_{\text{list}}^{\text{expr}_p}(\llbracket \alpha|_{\text{vpath}[v_r]} \rrbracket)) \in R_p(\tilde{D})$ .

For the “if” direction, note that  $(D, \alpha) \models \psi$  for all  $\psi \in \text{atoms}(Q') \cap \text{atoms}(\tilde{Q})$  holds since  $\text{vars}(\psi) \subseteq \text{dom}(\alpha)$  for all  $\psi \in \text{atoms}(Q') \cap \text{atoms}(\tilde{Q})$ . It remains to show the second part. Let  $R_p(u_1, \dots, u_t, x_{\text{expr}_p})$  be an arbitrary atom in  $\text{atoms}(\tilde{Q})$ . Since there is a  $\beta \supseteq \alpha$  such that  $(D, \beta) \models \psi$  for all  $\psi \in \text{atoms}(Q')$  it follows that  $\llbracket \gamma \rrbracket$  for all  $\gamma \subseteq \alpha$ , where  $\text{dom}(\gamma)$  forms a path in  $T'$ , is fit. It follows from Lemma 4.7 that for all  $u \in \text{child}(v^{\llbracket \gamma \rrbracket})$  that  $\mathcal{L}_u^{\llbracket \gamma \rrbracket}$  is not empty. Since the variable of  $\text{expr}_p$  is included in one of the sets  $\text{child}(v^{\llbracket \gamma \rrbracket})$  for one of those  $\gamma$ , it follows by definition of  $\tilde{D}$  that  $(\alpha(u_1), \dots, \alpha(u_r), w_{\text{list}}^{\text{expr}_p}(\llbracket \alpha|_{\text{vpath}[v_r]} \rrbracket)) \in R_p(\tilde{D})$ .

We consider now the “only if” direction. Let  $\tilde{T}$  be the induced subgraph of  $\tilde{T}$  on  $\{x_1, \dots, x_s\} = \text{dom}(\alpha)$ . We show by an induction over the height of the variables in  $\tilde{T}$  that the items  $\iota_v := \llbracket \alpha|_{\text{vpath}[v]} \rrbracket$  for all  $v \in V(T') \cap V(\tilde{T})$  are fit in the data structure

### 6.7. A reduction to q-hierarchical queries without aggregates

for  $Q'$  and  $D$ . Then, it follows that there is a  $\beta \supseteq \alpha$  such that  $(D, \beta) \models \psi$  for all  $\psi \in \text{atoms}(Q')$ .

For the induction base let us consider an item  $\iota_v$  of height 0. Then,  $v$  is either a leaf in  $\tilde{T}$  or the successor nodes of  $v$  are the nodes from  $\{x_1, \dots, x_q\}$ . Since  $v \in \text{vars}(\tilde{Q}) \cap \text{vars}(Q')$  we have  $\text{exatoms}(v) \subseteq \text{atoms}(Q') \cap \text{atoms}(\tilde{Q})$ . Since  $\text{exatoms}(v) \subseteq \text{atoms}(Q') \cap \text{atoms}(\tilde{Q})$ , it follows by assumption that for all  $\psi \in \text{exatoms}(v)$  holds  $(D, \alpha) \models \psi$ . For every child  $x_p \in \text{child}(v)$  ( $\text{child}(v)$  might be empty) it follows that  $x_p \in \{x_1, \dots, x_q\}$  and there is a tuple in  $R_p$  with  $(\alpha(u_1), \dots, \alpha(u_r), n_p) \in R_p^{\tilde{D}}$  where  $v_{\text{root}}, u_1, \dots, u_r$  is the path from  $v_{\text{root}}$  to  $x_p$  in  $T'$ . By definition, we have that  $\mathcal{L}_{x_p}^{[\alpha|_{\text{vpath}[u]}]} \neq \emptyset$ . It follows from Lemma 4.7 that  $\iota_v$  is fit.

For the inductive step let us consider an item  $\iota_v$  of height  $h$ . By the induction hypothesis we obtain that for all  $u \in \text{child}(v)$  the item  $\iota_u$  is fit and in particular, the lists  $\mathcal{L}_u^{\iota_v}$  are not empty. Since  $\text{exatoms}(v) \subseteq \text{atoms}(Q') \cap \text{atoms}(\tilde{Q})$ , it follows by assumption that for all  $\psi \in \text{exatoms}(v)$  holds  $(D, \alpha) \models \psi$ . It follows from Lemma 4.7 that  $\iota_v$  is fit.

**Proof of part (b).** It is straightforward to see that such a query can be easily computed in time  $O(\text{poly}(Q))$  after the q-tree of  $Q$  was computed.

**Proof of part (c).** To maintain  $Q$  on  $D$  under updates, we use two data structures in parallel. The data structure  $D_1$  is the data structure in Section 6.6 for  $Q$  on  $D$  and a data structure  $D_2$  for evaluation  $\tilde{Q}$  on  $\tilde{\sigma}$ -db  $\tilde{D}$  (with initialisation time  $t_i$  and update time  $t_u$ ).

When *initialising* the data structure for the empty database, we initialise  $D_1$  and  $D_2$  for the empty database. This takes time  $O(t_i)$ .

When the database  $D$  receives an update, we update the data structure  $D_1$ . If the update command contains a relation  $R \in \sigma \cap \tilde{\sigma}$ , we update  $D_2$ . During the update in  $D_1$ , we store for every item  $i$  where a  $w_{\text{list}}^{\text{expr}_p}(i)$ -value was changed, the old value  $\text{old}(w_{\text{list}}^{\text{expr}_p}(i))$  (or  $\text{nil}$  if it does not exist) and the new value  $\text{new}(w_{\text{list}}^{\text{expr}_p}(i))$  (or  $\text{nil}$  if it does not exist). Note that during an update step there are at most  $\text{poly}(Q)$  items, where the fit status is changing. In particular  $\text{poly}(Q)t_a$  aggregate values are changing in time  $\text{poly}(Q)t_a$ . For every change, we have to update the data structure for  $\tilde{Q}$  on  $\tilde{D}$  in time  $t_u$  by doing the updates  $\text{delete}R_p(\alpha^i(u_1), \dots, \alpha^i(u_q), \text{old}(w_{\text{list}}^{\text{expr}_p}(i)), b_1, \dots, b_p)$ , if  $\text{old}(w_{\text{list}}^{\text{expr}_p}(i)) \neq \text{nil}$  and  $\text{insert}R_p(\alpha^i(u_1), \dots, \alpha^i(u_q), \text{new}(w_{\text{list}}^{\text{expr}_p}(i)), b_1, \dots, b_p)$  if  $\text{new}(w_{\text{list}}^{\text{expr}_p}(i)) \neq \text{nil}$  for every item  $i$  where the  $w_{\text{list}}^{\text{expr}_p}(i)$  changes. All the update steps take time  $2\text{poly}(Q)t_u = \text{poly}(Q)t_u$ . In particular an update step takes time  $\text{poly}(Q)(t_a + t_u)$ .  $\square$

Now, we can give a proof for Theorem 3.7.

*Proof of Theorem 3.7.* Suppose that we receive as input a q-hierarchical CQ with aggregates  $Q$ . To initialise the data structure for  $Q$  on database  $D$ , the algorithm, we obtain from Lemma 6.18(b), computes a corresponding schema  $\tilde{\sigma}$  and a q-hierarchical query  $\tilde{Q}$  (without aggregates) and the q-tree using Lemma 6.8. Using Theorem 3.3, we obtain that there is a data structure for  $\tilde{Q}$  and  $\tilde{\sigma}$ -dbs that can be ini-

## 6. Answering $q$ -hierarchical conjunctive queries with aggregates under updates

tialised for the empty database in time  $\tilde{t}_i = O(\text{poly}(Q))$  and can be updated in time  $\tilde{t}_u = \text{poly}(\tilde{Q}) \leq \text{poly}(Q)$ . Using Lemma 6.18(c) we obtain an algorithm for  $Q$  with initialisation time  $t_i = O(\text{poly}(Q))$  and update time  $t_u = \text{poly}(Q)t_a$ . Note that the data structure from Theorem 3.3 for  $\tilde{Q}$  on  $\tilde{D}$  is present.

The **enumerate**-routine can be done as follows: Use the **enumerate**-routine from Theorem 3.3(b) for  $\tilde{Q}$  and  $\tilde{D}$ . Since  $Q(D) = \tilde{Q}(\tilde{D})$  (Lemma 6.18(a)), it follows that we enumerate the tuples in  $Q(D)$ .

The **test**-routine can be done as follows: Use the **test**-routine from Theorem 3.3(a) for  $\tilde{Q}$  and  $\tilde{D}$ . Since  $Q(D) = \tilde{Q}(\tilde{D})$  (Lemma 6.18(a)), it follows that we test upon input of a tuple  $\bar{a}$  if  $\bar{a} \in Q(D)$ .

The difference routine can be done as follows: Use the **diff**-routine from Theorem 3.3(c) for  $\tilde{Q}$  and  $\tilde{D}$ . Since  $Q(D^-) = \tilde{Q}(\tilde{D}^-)$  and  $Q(D^+) = \tilde{Q}(\tilde{D}^+)$  (Lemma 6.18(a)), it can enumerate the sets  $Q(D^+) \setminus Q(D^-)$  and  $Q(D^-) \setminus Q(D^+)$  with delay  $\text{poly}(Q)$ .

The **count**-routine can be done as follows: Use the **count**-routine from Theorem 3.3(a) for  $\tilde{Q}$  and  $\tilde{D}$ . Since  $Q(D) = \tilde{Q}(\tilde{D})$  (Lemma 6.18(a)), it follows that we test on input of a tuple  $\bar{a}$  if  $\bar{a} \in Q(D)$ .  $\square$

Using the counting expression in Section 6.5 and Theorem 3.7(a) we can give a proof for Theorem 3.3(d):

*Proof of Theorem 3.3(d).* On input of a  $Q$  and a database  $D$  compute the  $q$ -tree  $T_Q$  and the corresponding query  $Q_c$  with the counting expression. This can be computed in  $\text{poly}(Q)$ . Then use Theorem 3.7 to additionally maintain a data structure for  $Q_c$  and  $D$ . Clearly, since  $t_a = O(1)$ , we can maintain the data structure with initialisation time  $t_i = \text{poly}(Q)$  and  $t_u = \text{poly}(Q)$ . Additionally, we store a variable **value** where we store the result  $Q_c(D)$  after every update. This can be done by using the **enumerate** routine for  $Q_c$  and  $D$  in Theorem 3.7(a). The **count** routine can be done by outputting **value**. This can be done in  $O(1)$ .  $\square$

## 7. An application concerning learning polynomial regression models over q-hierarchical queries

In this chapter we show how to use aggregate queries from the previous chapter to efficiently prepare the training data (which is here the result set of a q-hierarchical conjunctive query) for learning a polynomial regression function, i.e., we prove Theorem 3.8 that is restated in the following.

**Theorem 3.8.** *For every  $d \in \mathbb{N}_{\geq 1}$  there is a dynamic algorithm that receives a q-hierarchical conjunctive query  $Q$  and a  $\sigma$ -db  $D_0$  and computes within  $t_i = O(\text{poly}(Q)^{2d+1})$  initialisation time and  $t_p = \text{poly}(Q)^{2d+1} \cdot O(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)^{2d+1}$  and allows to prepare the learning of the parameters of a polynomial hypothesis function  $h_\Theta$  of degree  $d$  in time  $t_l = O(\text{poly}(Q)^{2d+1})$  where  $D$  is the current database.*

We consider a problem, described by a function  $F : \mathbf{dom}^{k-1} \rightarrow \mathbf{dom}$  where we receive as input a  $(k-1)$ -ary tuple  $\bar{a}$  of  $\mathbf{dom}$  and map this to  $F(\bar{a})$ . In the regression task of machine learning we try to approximate the function  $F$  from a set

$$\left\{ (b^{(1)}, a_1^{(1)}, \dots, a_{k-1}^{(1)}), \dots, (b^{(m)}, a_1^{(m)}, \dots, a_{k-1}^{(m)}) \right\}.$$

This ultimate set is called training set and consists of tuples containing the input values  $a_1^{(i)}, \dots, a_{k-1}^{(i)}$  together with their results  $b^{(i)}$ . The  $a$ -values (the input) are called features and the  $b$ -value (the output) is called label. In our setting, the training set is taken from a query result  $Q(D)$  for a q-hierarchical conjunctive query

$$Q = \{(y, x_1, \dots, x_{k-1}, b_1, \dots, b_\ell) : \varphi\}$$

where  $\{y, x_1, \dots, x_{k-1}\} = \text{free}(\varphi)$  and  $b_1, \dots, b_\ell \in \mathbf{dom}$ . Note that the training set we use is not exactly the query result. We delete in every tuple the components that are constants. Using a feature or a label as constant, i.e., the value in every tuple in the result set is the same, forces that the trained model does not depend on the feature of the constant. This is the reason why we simply ignore constants in this context.

The goal is to define a function  $h_\Theta$  which is based on the training set and approximates a label for an arbitrary  $(k-1)$ -ary tuple of features. This is called *hypothesis function*.

Before we start, let us first consider an example that describes an application for machine learning and databases [85].

7. An application concerning learning polynomial regression models over q-hierarchical queries

**Example 7.1.** Let us assume that our aim is to consult companies, that have web pages, to get a better visibility in search engines by improving their positions. For most companies this is a very interesting question since most users often recognise the first few results of a search engine only. A feasible solution might be the following. We design a database  $D$  that consists of a relation  $R_{\text{Position}}$  of arity 4 and a relation  $R_{\text{textlength}}$  of arity 3.

Some software crawls the search engine on input of a keyword  $a_{\text{keyword}}$  and inserts for every position  $a_{\text{Pos}}$  the tuple  $(a_{\text{date}}, a_{\text{url}}, a_{\text{keyword}}, a_{\text{Pos}})$  into the relation  $R_{\text{Position}}(D)$  where  $a_{\text{url}}$  is the url of the website at position  $a_{\text{Pos}}$  when requesting the result in the search engine for the keyword  $a_{\text{keyword}}$  and  $a_{\text{date}}$  is the current date.

Another software analyses the websites and counts their number of words in the site. The software stores the result in the database by inserting the tuple  $(a_{\text{date}}, a_{\text{url}}, b_{\text{length}})$  to  $R_{\text{textlength}}(D)$  where  $b_{\text{length}}$  is the length of the web page, i.e., the number of words in the site, with url  $a_{\text{url}}$  and  $a_{\text{date}}$  is the current date.

Let us consider the following q-hierarchical conjunctive query

$$Q_{\text{SEO}} := \left\{ (y_{\text{length}}, x_{\text{date}}, x_{\text{url}}, x_{\text{Pos}}) : \exists x_{\text{keyword}} \left( \begin{array}{c} R_{\text{Position}} x_{\text{date}} x_{\text{url}} x_{\text{keyword}} x_{\text{Pos}} \wedge \\ R_{\text{textlength}} x_{\text{date}} x_{\text{url}} y_{\text{length}} \end{array} \right) \right\}.$$

The query result  $Q_{\text{SEO}}(D)$  is the set of tuples  $(b_{\text{length}}, a_{\text{date}}, a_{\text{url}}, a_{\text{Pos}})$  such that on  $a_{\text{date}}$  for the website with url  $a_{\text{url}}$  and length  $b_{\text{length}}$  there is a keyword such that  $a_{\text{url}}$  was in position  $a_{\text{Pos}}$  in the search engine on input of the keyword.

Taking the length as label and the other components as features, we can learn a function  $F$  where  $F(x_{\text{date}}, x_{\text{url}}, x_{\text{Pos}})$  approximates the length of the website  $x_{\text{url}}$  on position  $x_{\text{Pos}}$  on  $x_{\text{date}}$ . Moreover, we can use the function to check for another position, which text length we need, if we want to get a higher position in the future and we can consult the companies to change the length (In fact, there are a lot of other properties that may change the position in a search engine. For these properties, we can operate in the same manner as for the text length). Another useful application for this function is that we can take the first derivation of the function to receive correlations, such as a correlation between the position and the length.

Regression functions try to minimize the error between the value of the function and the training data. In particular, it might be possible that some tuples in the training data have too much influence in leaning the hypothesis function such that the polynomial does not give a good approximation for most of the training data, e.g. if the text length of the websites on position 1 to 9 and 11 to 20 is between 100 and 200 and the text length of the page in position 10 has length 5000000, the training data for position 10 will approximate the function  $F$  such that most of the values of  $F$  are over 200. That is the reason why the training data has to be smoothed, i.e. in our example we would delete the text length for the web page in position 10. To smooth the training data, we simply have to delete the corresponding tuples.

For convenience, we assume in this chapter that  $\mathbf{dom} = \mathbb{R}$  and we use for the aggregates **prod** and **sum** variants where  $\mathbf{dom} = \mathbb{R}$ , i.e., the binary relations  $\circ_{\text{sum}}$  and  $\circ_{\text{prod}}$  operate over real numbers.

We consider in Section 7.2 the case that the function  $h_\Theta$  is a polynomial function. As a warm up, we deal in Section 7.1 with linear hypothesis functions using examples.

## 7.1. Linear regression

The first part of this section describes the theoretical background of linear regression.

In the linear regression task the hypothesis function is a linear function of the form

$$h_\Theta(\bar{a}) := \theta_0 + \theta_1 a_1 + \cdots + \theta_{k-1} a_{k-1}.$$

To approximate the parameters  $\Theta := (\theta_0, \dots, \theta_{k-1})$  we minimize the least squares regression error function

$$\xi(\Theta) := \frac{1}{2} \sum_{i=1}^m \left( h_\Theta(\bar{a}^{(i)}) - y^{(i)} \right)^2 = \frac{1}{2} \sum_{i=1}^m \left( \theta_0 + \sum_{j=1}^{k-1} \theta_j a_j^{(i)} - b^{(i)} \right)^2.$$

Introducing new constants  $\theta_k := -1$  and  $a_0^{(i)} := 1$  and variables  $a_k^{(i)} := b^{(i)}$  for all  $i \in [m]$ , we can rewrite the error function in a compact way

$$\xi(\Theta) = \frac{1}{2} \sum_{i=1}^m \left( \sum_{j=0}^k \theta_j a_j^{(i)} \right)^2.$$

To compute the minimum of  $\xi(\Theta)$  we have to compute the jacobian matrix  $J(\xi)$  of  $\xi(\Theta)$ . For all  $\ell \in \{0, \dots, k-1\}$  the  $(1, (\ell+1))$ th element of  $J(\xi)$  is:

$$\begin{aligned} J(\xi)_{1, \ell+1} &= \frac{\delta}{\delta \theta_\ell} \frac{1}{2} \sum_{i=1}^m \left( \sum_{j=0}^k \theta_j a_j^{(i)} \right)^2 = \sum_{i=1}^m \left( \sum_{j=0}^k \theta_j a_j^{(i)} \right) a_\ell^{(i)} \\ &= \sum_{i=1}^m \sum_{j=0}^k \theta_j a_j^{(i)} a_\ell^{(i)} = \sum_{j=0}^k \sum_{i=1}^m \theta_j a_j^{(i)} a_\ell^{(i)} = \sum_{j=0}^k \theta_j \sum_{i=1}^m a_j^{(i)} a_\ell^{(i)}. \end{aligned}$$

In the algorithmic task of learning the parameters, we divide the algorithm into two passes. In the first pass, the *forward pass*, one has to compute the sums  $\sum_{i=1}^m a_j^{(i)} a_\ell^{(i)}$  and in the second pass, the *backward pass*, where one uses a numerical method to learn the parameters  $\Theta$ .

To learn the parameters  $\Theta$  we can use a numerical method such as batch gradient descent-method (BGD for short). In the backward pass, using the batch gradient descend method, we set for all  $\ell \in \{0, \dots, k-1\}$  the parameter  $\theta_\ell$  to  $\theta_\ell - \alpha J(\xi)_{1, \ell}$ . This will be repeated until convergence. The number  $\alpha$  is called the learning rate. It represents the step size of the gradient descent. If it is not small enough,  $\xi(\Theta)$  may not be smaller after an iteration and if  $\alpha$  is too small one need a lot iterations until convergence. The number  $\alpha$  may change after an iteration.

7. An application concerning learning polynomial regression models over  $q$ -hierarchical queries

We focus here on the forward pass. There, a naive algorithm has to compute the sum  $\sum_{i=1}^m a_j^{(i)} a_\ell^{(i)}$  for all  $\ell \in \{0, \dots, k-1\}$ . This takes time  $O(m)$ . Note that since the training set is taken from the query result, the number  $m$  is polynomial in the size of the database. We will show now that if the query is  $q$ -hierarchical, we can construct a data structure that can be updated in  $O(\text{poly}(Q))$  and allows to output the sum in time  $O(1)$ . For all  $u, w \in \text{vars}(Q)$  let

$$\begin{aligned} \text{Cofactor}[u, w] &:= \sum_{\beta \in \tilde{\mathcal{E}}^{\{0,1\}}} \beta(u) \cdot \beta(w) \\ \text{and } \text{Cofactor}[u, x_0] &:= \text{Cofactor}[x_0, u] := \sum_{\beta \in \tilde{\mathcal{E}}^{\{0,1\}}} \beta(u). \end{aligned}$$

Note that since the training set is  $Q(D)$  it holds that

$$\text{Cofactor}[x_j, x_\ell] = \sum_{i=1}^m a_j^{(i)} a_\ell^{(i)} \quad \text{for all } j, \ell \in \{0, \dots, k\}.$$

To compute the cofactor efficiently we use for all  $u, w \in \text{vars}(Q)$  the following  $q$ -hierarchical queries with aggregates:

$$Q_{u,w} := \{(\text{prod}(\text{lre}(u_1), \dots, \text{lre}(u_\ell)) : \varphi\}$$

and for all  $v \in \text{vars}(Q_{u,w})$  is  $\text{lre}(v)$  defined as follows: Let  $u_1, \dots, u_\ell$  be the children of  $v$  in  $T_{\text{free}}$  and for all  $i \in [\ell]$  is  $e_i := \text{lre}(u_i)$  if  $u_i \in \text{path}[u] \cup \text{path}[w]$  and  $e_i := c_{u_i}$  otherwise, where  $c_{u_i}$  is the aggregate expression for counting from Definition 6.10. For all  $m, \ell \in \mathbb{N}$  let  $\mathcal{A}_\ell^m = \{(h_\ell^m)_n\}_{n \in \mathbb{N}}$  be the aggregate with

$$(h_\ell^m)_n(\llbracket \langle a_{1,1}, a_{1,2}, \dots, a_{1,\ell} \rangle, \dots, \langle a_{n,1}, a_{n,2}, \dots, a_{n,\ell} \rangle \rrbracket) := \sum_{i=1}^n (a_{i,1})^m \cdot \prod_{k=2}^{\ell} a_{i,k}.$$

Note that the aggregate takes at most  $O(1)$  time.

- If  $v \in \{u, w\}$  we have  $\text{lre}(v) := \mathcal{A}_{\ell+1}^m(\langle v, e_1, \dots, e_\ell \rangle)$  where  $m = 2$  if  $u = w$  and  $m = 1$  otherwise and
- if  $v \notin \{u, w\}$  we have  $\text{sum}(\text{prod}(e_1, \dots, e_\ell))$ .

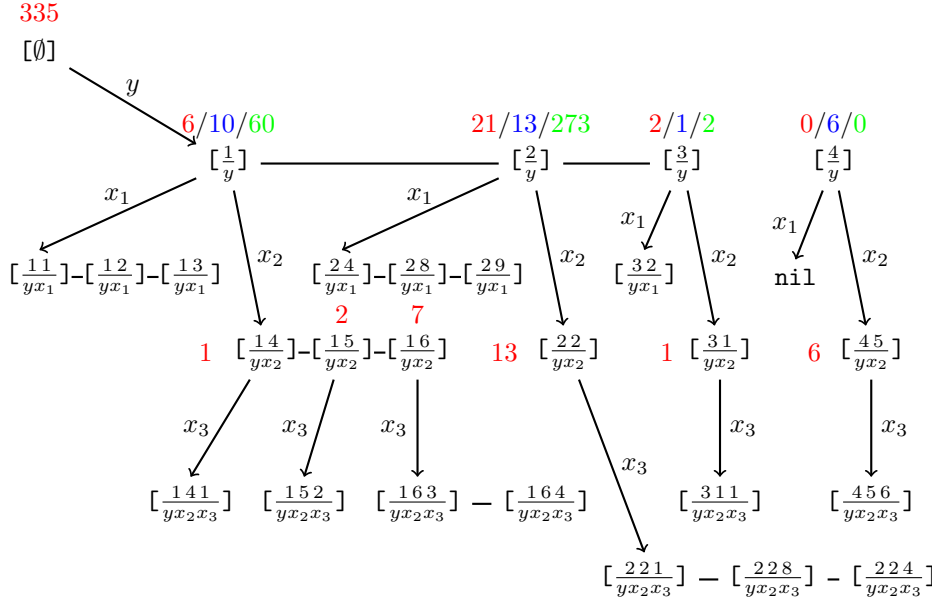
To show that for all  $x, y \in \text{vars}(Q)$  is  $Q_{x,y}(D) = \text{Cofactor}[x, y]$ , one can show that the following claim holds.

**Claim 7.2.** *Let  $Q$  be a  $q$ -hierarchical CQ and let  $D$  be a  $\sigma$ -db. For all items  $i$  in the data structure for  $Q$  on  $D$ , with  $v^i \in \text{path}[u] \cup \text{path}[w]$  holds for all  $p \in (\text{path}[u] \cup \text{path}[w]) \cap \text{child}(v^i)$  the following is satisfied.*

- If  $p \in \text{path}[u] \setminus \text{path}[w]$  then  $\llbracket \text{lre}(p) \rrbracket^{(D, \alpha^i)} = \sum_{\beta \in \tilde{\mathcal{E}}^i} \beta(u)$ .



Figure 7.1.: Illustration of  $Q_{x_1, x_3}$  on the database  $D$  of Example 4.4.



- If  $p \in \text{path}[w] \setminus \text{path}[u]$  then  $\llbracket \text{ire}(p) \rrbracket^{(D, \alpha^i)} = \sum_{\beta \in \tilde{\mathcal{E}}_i} \beta(w)$ .
- If  $p \in \text{path}[u] \cap \text{path}[w]$  then  $\llbracket \text{ire}(p) \rrbracket^{(D, \alpha^i)} = \sum_{\beta \in \tilde{\mathcal{E}}_i} \beta(u) \cdot \beta(w)$ .

Note that  $Q_{u,w}(D) = \text{Cofactor}[u, w]$  follows from Claim 7.2 for  $i = [\emptyset]$  since  $v_{\text{root}} \in \text{path}[u] \cap \text{path}[w]$ .

A proof of Claim 7.2 can be obtained by the proof of Lemma 7.3 for polynomial function by considering the special case that  $d = 1$ .

To get an idea how it works let us consider the database and the query from Example 4.4. For the example, we are interested in  $\text{Cofactor}[x_1, x_3]$  and  $\text{Cofactor}[x_3, x_3]$ . The corresponding queries are the following.

$$\begin{aligned} Q_{x_1, x_3} &= \{(\text{prod}(\text{sum}(\text{prod}(\mathcal{A}_1^1(\langle x_1 \rangle), \text{sum}(\text{prod}(\mathcal{A}_1^1(\langle x_3 \rangle)))))) : \varphi\} \\ Q_{x_3, x_3} &= \{(\text{prod}(\text{sum}(\text{prod}(\text{count}(x_1), \text{sum}(\text{prod}(\mathcal{A}_1^1(\langle x_3 \rangle)))))) : \varphi\} \\ \text{where } \varphi &:= Eyx_1 \wedge Fyx_2x_3 \wedge Gyx_2x_3 \end{aligned}$$

See Figure 7.1 for an illustration of  $Q_{x_1, x_3}$  on the database  $D_0$  of Example 4.4. A red number on the left or above an item  $i$  with  $v^i = x_2$  represents the number  $w_{\text{list}}^{\mathcal{A}_1^1(\langle x_3 \rangle)}(i)$ , i.e., the number of the sum over  $\iota \in \mathcal{L}_{x_3}^i$  of their constants  $a^\iota$ . Since  $x_2$  has only one child it follows that  $w_{\text{item}}^{\text{prod}(\mathcal{A}_1^1(\langle x_3 \rangle))}(i) = w_{\text{list}}^{\mathcal{A}_1^1(\langle x_3 \rangle)}(i)$  for all items  $i$  with  $v^i = x_2$ . The three numbers  $n_1/n_2/n_3$  above an item  $i$  with  $v^i = y$  represents that  $n_1 =$

7. An application concerning learning polynomial regression models over  $q$ -hierarchical queries

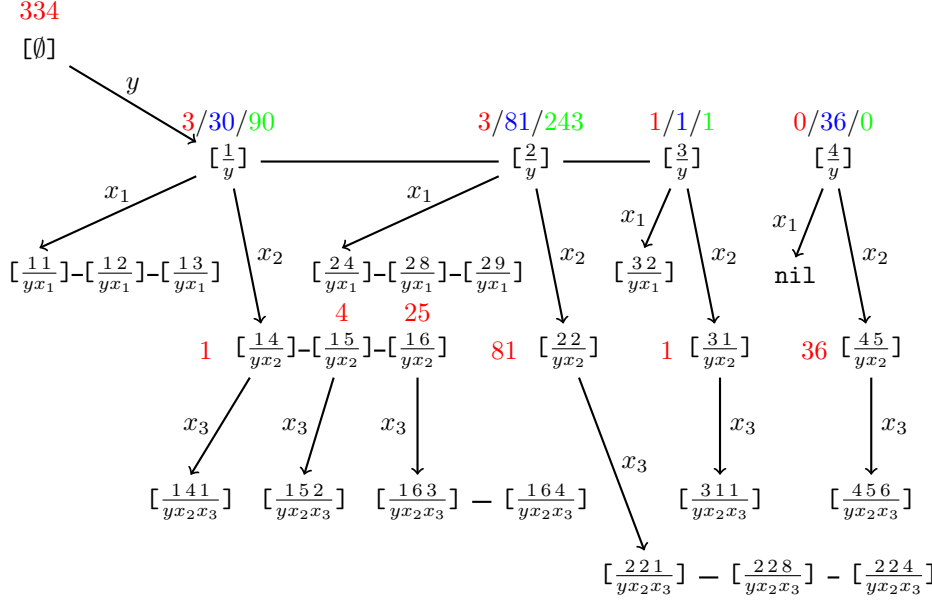
$y$	$x_1$	$x_2$	$x_3$	$x_1 \cdot x_3$	$(x_3)^2$
1	1	4	1	1	1
1	1	5	2	2	4
1	1	6	3	3	9
1	1	6	4	4	16
1	2	4	1	2	1
1	2	5	2	4	4
1	2	6	3	6	9
1	2	6	4	8	16
1	3	4	1	3	1
1	3	5	2	6	4
1	3	6	3	9	9
1	3	6	4	12	16
2	4	2	1	4	1
2	4	2	8	32	64
2	4	2	4	16	16
2	8	2	1	8	1
2	8	2	8	64	64
2	8	2	4	32	16
2	9	2	1	9	1
2	9	2	8	72	64
2	9	2	4	36	16
3	2	3	1	2	1
				$\Sigma =$	$\Sigma =$
				<b>335</b>	<b>334</b>

Table 7.1.: Enumeration of  $Q(D_0)$  from Example 4.4 together with the product  $x_1 \cdot x_3$  and  $x_3 \cdot x_3$  for every result tuple and the sum over all result tuples of these products.

$w_{\text{list}}^{\mathcal{A}_1^1(\langle x_1 \rangle)}(i)$ , i.e.,  $n_1$  is the number of the sum of  $a^\iota$  over  $\iota \in \mathcal{L}_{x_1}^i$ , the number  $n_2 = w_{\text{list}}^{\text{sum}(\text{prod}(\mathcal{A}_1^1(\langle x_3 \rangle)))}(i)$ , i.e., the sum over  $\iota \in \mathcal{L}_{x_2}^i$  of their  $w_{\text{item}}^{\text{prod}(\mathcal{A}_1^1(\langle x_3 \rangle))}(\iota)$  values. The number  $n_3$  is  $w_{\text{list}}^{\text{prod}(\mathcal{A}_1^1(\langle x_1 \rangle), \text{sum}(\text{prod}(\mathcal{A}_1^1(\langle x_3 \rangle))))}(i)$ , i.e., the product of  $n_1$  and  $n_2$ . The red number above the start item depicts  $w_{\text{list}}^{\text{sum}(\text{prod}(\mathcal{A}_1^1(\langle x_1 \rangle), \text{sum}(\text{prod}(\mathcal{A}_1^1(\langle x_3 \rangle))))}([\emptyset])$ . Note that since the node  $v_{\text{root}}$  has only one child in the  $q$ -tree of  $Q_{x_1, x_2}$ , the number is also the query result  $Q_{x_1, x_2}(D)$ .

See Table 7.1 for the query result of  $Q(D)$  and the product  $\beta(x_1) \cdot \beta(x_3)$  and  $(\beta(x_3))^2$  for every tuple in the result set. These products are given in the line of the corresponding tuple. At the bottom of the columns of  $x_1 \cdot x_3$  and  $(x_3)^2$  the sum over all products is given, which is equal to  $\sum_{\beta \in \tilde{\mathcal{E}}_{\{0\}}} \beta(x_1) \cdot \beta(x_3)$  (resp.  $\sum_{\beta \in \tilde{\mathcal{E}}_{\{0\}}} \beta(x_3)^2$ )

Let us consider the item  $i = [\frac{1}{y}]$ . Recall the value  $n_3$  of  $i$  from the previous

Figure 7.2.: Illustration of  $Q_{x_3, x_3}$  on the database  $D$  of Example 4.4.

description:

$$\begin{aligned}
 n_3 &= (1 + 2 + 3)(1 + 2 + (3 + 4)) \\
 &= 1 \cdot 1 + 1 \cdot 2 + 1 \cdot 3 + 1 \cdot 4 + 2 \cdot 1 + 2 \cdot 2 + 2 \cdot 3 + 2 \cdot 4 + 3 \cdot 1 + 3 \cdot 2 + 3 \cdot 3 + 3 \cdot 4
 \end{aligned}$$

Compare these values with those in Table 7.1, we see, that  $n_3$  is the sum over all assignments  $\beta \in \tilde{\mathcal{E}}^{[\emptyset]}$  with  $\beta(y) = 1$  of  $\beta(x_1) \cdot \beta(x_3)$ . The number  $n_3$  is equal to  $\sum_{\beta \in \tilde{\mathcal{E}}^{[\frac{1}{y}]}} \beta(x_1) \cdot \beta(x_3)$  since the decomposition lemma (Lemma 4.12) states that each assignment in  $\tilde{\mathcal{E}}^{[\frac{1}{y}]}$  is the union of assignments  $\beta_1$  and  $\beta_2$  where for  $\beta_1$  there is an item  $\iota_1 \in \mathcal{L}_{x_1}^{[\frac{1}{y}]}$  such that  $\beta_1 \in \tilde{\mathcal{E}}^{\iota_1}$  and for  $\beta_2$  there is an item  $\iota_2 \in \mathcal{L}_{x_2}^{[\frac{1}{y}]}$  such that  $\beta_2 \in \tilde{\mathcal{E}}^{\iota_2}$  and thus  $\sum_{\beta \in \tilde{\mathcal{E}}^{[\frac{1}{y}]}} \beta(x_1) \cdot \beta(x_3) = \sum_{\iota_1 \in \mathcal{L}_{x_1}^{[\frac{1}{y}]}} \sum_{\beta_1 \in \tilde{\mathcal{E}}^{\iota_1}} \beta_1(x_1) \sum_{\iota_2 \in \mathcal{L}_{x_2}^{[\frac{1}{y}]}} \sum_{\beta_2 \in \tilde{\mathcal{E}}^{\iota_2}} \beta_2(x_2) = n_1 \cdot n_2$ .

The same can be shown for the items  $[\frac{2}{y}]$  and  $[\frac{3}{y}]$ . Since

$$w_{\text{list}}^{\text{sum}(\text{prod}(\mathcal{A}_1^1(\langle x_1 \rangle)), \text{sum}(\text{prod}(\mathcal{A}_1^1(\langle x_3 \rangle))))}([\emptyset])$$

is the sum of all the  $n_3$  values of  $[\frac{1}{y}]$ ,  $[\frac{2}{y}]$  and  $[\frac{3}{y}]$  it follows that the query result is  $\sum_{\beta \in \tilde{\mathcal{E}}^{[\emptyset]}} \beta(x_1) \cdot \beta(x_3)$ .

Let us now consider the query  $Q_{x_3, x_3}$  and the corresponding Figure 7.2.

A red number on the left or above an item  $i$  with  $v^i = x_2$  depicts the number  $w_{\text{list}}^{\mathcal{A}_1^2(\langle x_3 \rangle)}(i)$ , i.e., the number of the sum over  $\iota \in \mathcal{L}_{x_3}^i$  of the square of their constants

## 7. An application concerning learning polynomial regression models over $q$ -hierarchical queries

$a^\iota$ . Note that since  $x_2$  has only one child it follows that  $w_{\text{item}}^{\text{prod}(\mathcal{A}_1^2(\langle x_3 \rangle))}(i) = w_{\text{list}}^{\mathcal{A}_1^2(\langle x_3 \rangle)}(i)$  for all items  $i$  with  $v^i = x_2$ . The three numbers  $n_1/n_2/n_3$  above an item  $i$  with  $v^i = y$  depict that  $n_1 = w_{\text{list}}^{\text{count}(x_1)}(i)$ , i.e., is the number of item in the  $\mathcal{L}_{x_1}^i$ -list, the number  $n_2 = w_{\text{list}}^{\text{sum}(\text{prod}(\mathcal{A}_1^2(\langle x_3 \rangle)))}(i)$ , i.e., the sum over  $\iota \in \mathcal{L}_{x_2}^i$  of their  $w_{\text{item}}^{\text{prod}(\mathcal{A}_1^2(\langle x_3 \rangle))}(\iota)$  values. The number  $n_3$  is  $w_{\text{list}}^{\text{prod}(\text{count}(x_1), \text{sum}(\text{prod}(\mathcal{A}_1^2(\langle x_3 \rangle))))}(i)$ , i.e., the product of  $n_1$  and  $n_2$ . The red number above the start item depicts  $w_{\text{list}}^{\text{sum}(\text{prod}(\text{count}(x_1), \text{sum}(\text{prod}(\mathcal{A}_1^2(\langle x_3 \rangle))))}([\emptyset])$ . Note that since the node  $v_{\text{root}}$  has only one child in the  $q$ -tree of  $Q_{x_3, x_3}$ , the number is equivalent to the query result  $Q_{x_3, x_3}(D)$ .

Let us consider the item  $i = [\frac{1}{y}]$ . Recall that the  $n_3$ -value of  $i$  is computed the following way:

$$n_3 := 3 \cdot (1 + 4 + (9 + 16))$$

Comparing this with Table 7.1 we see, that the values  $(x_3)^2$  will be repeated 3 times. This depends on the fact that the values for  $x_3$  will be repeated since every item in all assignments  $\beta_1 \in \tilde{\mathcal{E}}^{\iota_1}$  for every  $\iota_1 \in \mathcal{L}_{x_1}^i$  has to be compared with  $\beta_2 \in \tilde{\mathcal{E}}^{\iota_2}$  for every  $\iota_2 \in \mathcal{L}_{x_2}^i$  to receive an assignment in  $\tilde{\mathcal{E}}^i$ . In particular, every  $\beta_2 \in \tilde{\mathcal{E}}^{\iota_2}$  for every  $\iota_2 \in \mathcal{L}_{x_2}^i$  will be enumerated three times since there are three assignments in  $\left| \bigcup_{\iota_1 \in \mathcal{L}_{x_1}^i} \tilde{\mathcal{E}}^{\iota_1} \right|$ . This is the number given by  $n_1$ . Analogously, this can be shown for the items  $[\frac{2}{y}]$  and  $[\frac{3}{y}]$  and it is straightforward to see that their sums are equal to  $\sum_{\beta \in \tilde{\mathcal{E}}^{\{\emptyset\}}} \beta(x_3)^2$ .

## 7.2. Polynomial regression

In this section we enrich the idea for maintaining linear regression to polynomial regression, i.e., we show the result for polynomial regression of degree  $d$ . The hypothesis function for a polynomial of degree  $d$  is the following.

$$h_{\Theta}(a_1, \dots, a_k) = \sum_{j_1=0}^k \sum_{j_2=j_1}^k \dots \sum_{j_d=j_{d-1}}^k \theta_{(j_1, \dots, j_d)} \cdot a_{j_1} \cdot \dots \cdot a_{j_d}$$

where  $a_0 := 1$  and  $\Theta = (\theta_i)_{i \in \mathcal{I}_d}$  such that  $\mathcal{I}_d$  denotes the index set for  $d$  which is defined as  $\mathcal{I}_d := \{(j_1, \dots, j_d) : j_1, \dots, j_d \in \{0, \dots, k\}, j_1 \leq \dots \leq j_d\}$ . Using the least squares error function we receive

$$\begin{aligned} \xi(\Theta) &:= \frac{1}{2} \sum_{i=1}^m \left( h_{\Theta}(\bar{a}^{(i)}) - b^{(i)} \right)^2 \\ &= \frac{1}{2} \sum_{i=1}^m \left( \sum_{(j_1, \dots, j_d) \in \mathcal{I}_d} \theta_{(j_1, \dots, j_d)} \cdot a_{j_1}^{(i)} \cdot \dots \cdot a_{j_d}^{(i)} - b^{(i)} \right)^2 \end{aligned}$$

To compute the minimum of the error function, we consider the Jacobi matrix  $J(\xi)$  of  $\xi(\Theta)$ . The  $(1, (\ell_1, \dots, \ell_d))$ th component of the Jacobi matrix is

$$\begin{aligned}
J(\xi)_{1,(\ell_1, \dots, \ell_d)} &= \sum_{i=1}^m \left( \sum_{(j_1, \dots, j_d) \in \mathcal{I}_d} \theta_{(j_1, \dots, j_d)} \cdot a_{j_1}^{(i)} \cdot \dots \cdot a_{j_d}^{(i)} - b^{(i)} \right) \cdot a_{\ell_1}^{(i)} \cdot \dots \cdot a_{\ell_d}^{(i)} \\
&= \sum_{i=1}^m \sum_{(j_1, \dots, j_d) \in \mathcal{I}_d} \theta_{(j_1, \dots, j_d)} \cdot a_{j_1}^{(i)} \cdot \dots \cdot a_{j_d}^{(i)} \cdot a_{\ell_1}^{(i)} \cdot \dots \cdot a_{\ell_d}^{(i)} \\
&\quad - \sum_{i=1}^m b^{(i)} \cdot a_{\ell_1}^{(i)} \cdot \dots \cdot a_{\ell_d}^{(i)} \\
&= \sum_{(j_1, \dots, j_d) \in \mathcal{I}_d} \theta_{(j_1, \dots, j_d)} \sum_{i=1}^m a_{j_1}^{(i)} \cdot \dots \cdot a_{j_d}^{(i)} \cdot a_{\ell_1}^{(i)} \cdot \dots \cdot a_{\ell_d}^{(i)} \\
&\quad - \sum_{i=1}^m b^{(i)} \cdot a_{\ell_1}^{(i)} \cdot \dots \cdot a_{\ell_d}^{(i)}
\end{aligned}$$

To learn the parameters in the backward pass, one can use the batch gradient descent method (see Algorithm 10).

---

**Algorithm 10** Learn the parameters  $\Theta$  using the batch gradient descent method on polynomial regression.

---

```

1: repeat
2:   for all  $(\ell_1, \dots, \ell_d) \in \mathcal{I}_d$  do
3:     Set  $\theta_{(\ell_1, \dots, \ell_d)}$  to  $\theta_{(\ell_1, \dots, \ell_d)} - \alpha J(\xi)_{1,(\ell_1, \dots, \ell_d)}$ 
4: until convergence

```

---

The number  $\alpha$  is once again the learning rate. We discuss how to efficiently maintain the sums of the form  $\sum_{i=1}^m a_{j_1}^{(i)} \cdot \dots \cdot a_{j_d}^{(i)} \cdot a_{\ell_1}^{(i)} \cdot \dots \cdot a_{\ell_d}^{(i)}$  under updates. Since the training set is taken from the query result, we obtain that

$$\sum_{i=1}^m a_{j_1}^{(i)} \cdot \dots \cdot a_{j_d}^{(i)} \cdot a_{\ell_1}^{(i)} \cdot \dots \cdot a_{\ell_d}^{(i)} = \sum_{\beta \in \tilde{\mathcal{E}}^{[\emptyset]}} \prod_{p=1}^d \beta(x_{j_p}) \cdot \beta(x_{\ell_p})$$

and

$$\sum_{i=1}^m b^{(i)} \cdot a_{\ell_1}^{(i)} \cdot \dots \cdot a_{\ell_d}^{(i)} = \sum_{\beta \in \tilde{\mathcal{E}}^{[\emptyset]}} \beta(y) \prod_{p=1}^d \beta(x_{\ell_p})$$

Here, we set  $\beta(x_0) := 1$  for all  $\beta \in \tilde{\mathcal{E}}^{[\emptyset]}$ . The following lemma shows that these sums can be maintained under updates.

7. An application concerning learning polynomial regression models over  $q$ -hierarchical queries

**Lemma 7.3.** *There is a dynamic algorithm that receives a  $q$ -hierarchical CQ  $Q$  and a  $\sigma$ -db  $D_0$  and computes within  $t_i = \text{poly}(Q)^{2d}$  initialisation time, and  $t_p = \text{poly}(Q)^{2d} \|D\|$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)$  and allows upon input of a tuple  $(w_1, \dots, w_{\tilde{d}}) \in \text{vars}(Q)^{\tilde{d}}$  for  $\tilde{d} \leq 2d$  the value*

$$\sum_{\beta \in \tilde{\mathcal{E}}^{\{\emptyset\}}} \prod_{p=1}^{\tilde{d}} \beta(w_p)$$

in time  $O(1)$ .

*Proof.* To handle these sums under updates, we define for every  $\tilde{d} \leq 2d$  and every  $(w_1, \dots, w_{\tilde{d}}) \in \text{vars}(Q)^{\tilde{d}}$  a query  $Q_{(w_1, \dots, w_{\tilde{d}})}$  such that

$$Q_{(w_1, \dots, w_{\tilde{d}})}(D) = \sum_{\beta \in \tilde{\mathcal{E}}^{\{\emptyset\}}} \prod_{p=1}^{\tilde{d}} \beta(w_p).$$

In the following we show now how to define a query  $Q_{\bar{w}}$  for an arbitrary tuple  $\bar{w} = (w_1, \dots, w_{\tilde{d}})$ . Let  $T_{\text{free}}$  be the subgraph of the  $q$ -tree of  $Q$  induced on  $\text{free}(Q)$  and let  $\tilde{T}_{\bar{w}}$  be the subgraph of  $T_{\text{free}}$  induced on  $\bigcup_{i=1}^{\tilde{d}} \text{path}[w_i]$ . For all  $v \in \text{vars}(Q)$  we define  $\#v$  as the number of appearance of  $v$  in the tuple  $\bar{w}$ .

**Definition 7.4.** *The aggregate expression  $\text{reg}_{\bar{w}}$  is defined as follows. For all  $v \in V(\tilde{T}_{\bar{w}})$  let  $u_1, \dots, u_s$  be the children of  $v$  in  $T_{\text{free}}$  and  $e_i := \text{reg}_{\bar{w}}(u_i)$  if  $u_i \in V(\tilde{T}_{\bar{w}})$  and  $e_i := c_{u_i}$  otherwise where  $c_{u_i}$  is the counting expression from Definition 6.10. Then,*

- if  $v \in \{w_1, \dots, w_{\tilde{d}}\}$  then  $\text{reg}_{\bar{w}}(v) := \mathcal{A}_{s+1}^{\#v}(\langle v, e_1, \dots, e_s \rangle)$  and
- if  $v \notin \{w_1, \dots, w_{\tilde{d}}\}$  then  $\text{reg}_{\bar{w}}(v) := \text{sum}(\text{prod}(e_1, \dots, e_s))$ .

For any  $q$ -hierarchical conjunctive query  $Q = \{(x_1, \dots, x_k) : \varphi\}$  the query  $Q_{\bar{w}} = \{(\text{prod}(e_1, \dots, e_s)) : \varphi\}$  where  $z_1, \dots, z_s$  are the children of the root in the  $q$ -tree of  $Q$  and  $e_j := \text{reg}_{\bar{w}}(z_j)$  if  $z_j \in V(\tilde{T}_{\bar{w}})$  and  $e_j := c_{z_j}$  otherwise.

Our aim is to show that

$$Q_{(w_1, \dots, w_{\tilde{d}})}(D) = \sum_{\beta \in \tilde{\mathcal{E}}^{\{\emptyset\}}} \prod_{p=1}^{\tilde{d}} \beta(w_p).$$

The following notation and the following claim will be convenient for the proof.

For all  $v \in V(\tilde{T}_{\bar{w}})$  let

$$M_v := \text{succ}(v) \cap V(\tilde{T}_{\bar{w}}) \quad \text{and} \quad W_v := M_v \cap \{w_1, \dots, w_{\tilde{d}}\}.$$

**Claim 7.5.** For all items  $i$  with  $v^i \in V(\tilde{T}_{\bar{w}})$  we have

$$\llbracket \text{reg}_{\bar{w}}(u) \rrbracket^{(D, \alpha^i)} = \sum_{\iota \in \mathcal{L}_u^i} \sum_{\beta \in \tilde{\mathcal{E}}^\iota} \prod_{z \in W_u} \beta(z)^{\#z}$$

for all  $u \in M_{v^i} \cap \text{child}(v^i)$ .

We show this claim by induction of the height of an item in  $\tilde{T}_{\bar{w}}$ .

For the induction base let us consider an item  $i$  where  $v^i$  has height 1 in  $\tilde{T}_{\bar{w}}$ . Then, for all  $u \in M_{v^i}$  it holds that  $u \in \{w_1, \dots, w_d\}$ . Let  $z_1, \dots, z_s$  where  $s \geq 0$  be the children of  $v^i$  in  $T_{\text{free}}$ . Note that  $e_j = c_{z_j}$  for all  $j \in [s]$  since  $z_j \notin V(\tilde{T}_{\bar{w}})$ .

$$\begin{aligned} \llbracket \text{reg}_{\bar{w}}(u) \rrbracket^{(D, \alpha^i)} &= \sum_{\iota \in \mathcal{L}_u^i} (a^\iota)^{\#u} \cdot \prod_{j \in [s]} \llbracket c_{v_j} \rrbracket^{(D, \alpha^\iota)} \\ &\stackrel{(*)}{=} \sum_{\iota \in \mathcal{L}_u^i} (a^\iota)^{\#u} \cdot \prod_{j \in [s]} \left| \bigcup_{i \in \mathcal{L}_{v_j}^\iota} \tilde{\mathcal{E}}^i \right| \stackrel{(**)}{=} \sum_{\iota \in \mathcal{L}_u^i} (a^\iota)^{\#u} \cdot |\tilde{\mathcal{E}}^\iota| \\ &\stackrel{(***)}{=} \sum_{\iota \in \mathcal{L}_u^i} \sum_{\beta \in \tilde{\mathcal{E}}^\iota} \beta(u)^{\#u} \stackrel{(***)}{=} \sum_{\iota \in \mathcal{L}_u^i} \sum_{\beta \in \tilde{\mathcal{E}}^\iota} \prod_{z \in W_u} \beta(z)^{\#z} \end{aligned}$$

Equation  $(*)$  follows from Claim 6.13, and Equation  $(**)$  from Claim 6.12, and Equation  $(***)$  from the fact that  $|\tilde{\mathcal{E}}^\iota| = \sum_{\beta \in \tilde{\mathcal{E}}^\iota} 1$  and for all  $\iota \in \mathcal{L}_u^i$  and for all  $\beta \in \tilde{\mathcal{E}}^\iota$  we have  $\beta(u) = \alpha^\iota$ . Equation  $(***)$  follows from the fact that  $W_u = \{u\}$ .

For the inductive step we consider an item  $i$  where  $v^i$  has height  $h > 1$  in  $\tilde{T}_{\bar{w}}$ . Let  $u \in M_{v^i}$  be arbitrary and let  $v_1, \dots, v_s$  be the children of  $u$  in  $T$ . Let us consider the case that  $u \in \{w_1, \dots, w_d\}$ . Then, it follows that

$$\begin{aligned} \llbracket \text{reg}_{\bar{w}}(u) \rrbracket^{(D, \alpha^i)} &= \sum_{\iota \in \mathcal{L}_u^i} (a^\iota)^{\#u} \cdot \prod_{\substack{j \in [s] \\ v_j \notin M_u}} \llbracket c_{v_j} \rrbracket^{(D, \alpha^\iota)} \cdot \prod_{\substack{j \in [s] \\ v_j \in M_u}} \llbracket \text{reg}_{\bar{w}}(v_j) \rrbracket^{(D, \alpha^\iota)} \\ &\stackrel{(*)}{=} \sum_{\iota \in \mathcal{L}_u^i} (a^\iota)^{\#u} \cdot \prod_{\substack{j \in [s] \\ v_j \notin M_u}} \left| \bigcup_{i \in \mathcal{L}_{v_j}^\iota} \tilde{\mathcal{E}}^i \right| \cdot \prod_{\substack{j \in [s] \\ v_j \in M_u}} \sum_{i \in \mathcal{L}_{v_j}^\iota} \sum_{\gamma \in \tilde{\mathcal{E}}^i} \prod_{z \in W_{v_j}} \gamma(z)^{\#z} \\ &\stackrel{(**)}{=} \sum_{\iota \in \mathcal{L}_u^i} \sum_{\beta \in \tilde{\mathcal{E}}^\iota} \prod_{z \in W_u} \beta(z)^{\#z} \end{aligned}$$

Equation  $(*)$  follows from Claim 6.13 and the induction hypothesis.

Equation  $(**)$  follows from the following fact. From the decomposition lemma (Lemma 4.12) we know that for every  $\iota \in \mathcal{L}_u^i$  every assignment  $\beta \in \tilde{\mathcal{E}}^\iota$  can be decomposed into  $\beta = \beta_u \cup \beta_1 \cup \dots \cup \beta_s$  such that there is a  $\iota_j \in \mathcal{L}_{v_j}^\iota$  where  $\beta_j \in \tilde{\mathcal{E}}^{\iota_j}$  for

7. An application concerning learning polynomial regression models over  $q$ -hierarchical queries

all  $j \in [s]$ . Therefore, it follows

$$\begin{aligned}
\sum_{\iota \in \mathcal{L}_u^i} \sum_{\beta \in \tilde{\mathcal{E}}^\iota} \prod_{z \in W_u} \beta(z)^{\#z} &\stackrel{4.12}{=} \sum_{\iota \in \mathcal{L}_u^i} \sum_{\iota_1 \in \mathcal{L}_{v_1}^\iota} \sum_{\beta_1 \in \tilde{\mathcal{E}}^{\iota_1}} \cdots \sum_{\iota_s \in \mathcal{L}_{v_s}^\iota} \sum_{\beta_s \in \tilde{\mathcal{E}}^{\iota_s}} \prod_{z \in W_u} \beta(z)^{\#z} \\
&\quad \text{where } \beta := \alpha^\iota \cup \beta_1 \cup \cdots \cup \beta_s \\
&= \sum_{\iota \in \mathcal{L}_u^i} \beta(u)^{\#u} \sum_{\iota_1 \in \mathcal{L}_{v_1}^\iota} \sum_{\beta_1 \in \tilde{\mathcal{E}}^{\iota_1}} \prod_{z \in W_u \cap \text{succ}(v_1)} \beta(z)^{\#z} \cdots \\
&\quad \sum_{\iota_s \in \mathcal{L}_{v_s}^\iota} \sum_{\beta_s \in \tilde{\mathcal{E}}^{\iota_s}} \prod_{z \in W_u \cap \text{succ}(v_s)} \beta(z)^{\#z} \\
&= \sum_{\iota \in \mathcal{L}_u^i} (a^\iota)^{\#u} \prod_{\substack{j \in [s] \\ v_j \notin M_u}} \sum_{\iota_j \in \mathcal{L}_{v_j}^\iota} \sum_{\beta_j \in \tilde{\mathcal{E}}^{\iota_j}} \prod_{z \in W_{v_j}} \beta_j(z)^{\#z} \cdot \\
&\quad \prod_{\substack{j \in [s] \\ v_j \in M_u}} \sum_{\iota_s \in \mathcal{L}_{v_s}^\iota} \sum_{\beta_s \in \tilde{\mathcal{E}}^{\iota_s}} \prod_{z \in W_{v_j}} \beta_j(z)^{\#z}
\end{aligned} \tag{7.1}$$

The last equation holds since  $\beta_j \subseteq \beta$  and  $\text{dom}(\beta_j) \supseteq W_{v_j}$  and  $W_{v_j} = (W_u \cap \text{succ}(v_j))$ . For all  $j \in [s]$  with  $v_j \notin M_u$  it follows that the product  $\prod_{z \in W_u \cap \text{succ}(v_j)} \beta(z)^{\#z}$  is empty since  $v_j \notin M_u$  and thus  $W_u \cap \text{succ}(v_j) = \emptyset$  and therefore we obtain

$$\begin{aligned}
\sum_{\iota_j \in \mathcal{L}_{v_j}^\iota} \sum_{\beta_j \in \tilde{\mathcal{E}}^{\iota_j}} \prod_{z \in W_{v_j}} \beta_j(z)^{\#z} &= \sum_{\iota_j \in \mathcal{L}_{v_j}^\iota} \sum_{\beta_j \in \tilde{\mathcal{E}}^{\iota_j}} 1 = \sum_{\iota_j \in \mathcal{L}_{v_j}^\iota} |\tilde{\mathcal{E}}^{\iota_j}| \\
&= \left| \bigcup_{\iota_j \in \mathcal{L}_{v_j}^\iota} \tilde{\mathcal{E}}^{\iota_j} \right|
\end{aligned} \tag{7.2}$$

The last equation follows from the fact that  $\tilde{\mathcal{E}}^{\iota_1} \cap \tilde{\mathcal{E}}^{\iota_2} = \emptyset$  for  $\iota_1, \iota_2 \in \mathcal{L}_u^i$  with  $\iota_1 \neq \iota_2$ . Comparing equations (7.1) and (7.2) we obtain that

$$\begin{aligned}
&\sum_{\iota \in \mathcal{L}_u^i} \sum_{\beta \in \tilde{\mathcal{E}}^\iota} \prod_{z \in W_u} \beta(z)^{\#z} \\
&= \sum_{\iota \in \mathcal{L}_u^i} (a^\iota)^{\#u} \prod_{\substack{j \in [s] \\ v_j \notin M_u}} \left| \bigcup_{\iota_j \in \mathcal{L}_{v_j}^\iota} \tilde{\mathcal{E}}^{\iota_j} \right| \cdot \prod_{\substack{j \in [s] \\ v_j \in M_u}} \sum_{\iota_s \in \mathcal{L}_{v_s}^\iota} \sum_{\beta_s \in \tilde{\mathcal{E}}^{\iota_s}} \prod_{z \in W_{v_j}} \beta_j(z)^{\#z}
\end{aligned}$$

This concludes the proof of the correctness of  $(\star\star)$ . Let us now consider the last case that  $u \notin \{w_1, \dots, w_d\}$ . Then,

$$\llbracket \text{reg}_{\bar{w}}(u) \rrbracket^{(D, \alpha^i)} = \sum_{\iota \in \mathcal{L}_u^i} \prod_{\substack{j \in [s] \\ v_j \notin M_u}} \llbracket c_{v_j} \rrbracket^{(D, \alpha^\iota)} \cdot \prod_{\substack{j \in [s] \\ v_j \in M_u}} \llbracket \text{reg}_{\bar{w}}(v_j) \rrbracket^{(D, \alpha^\iota)}$$



The proof for that case can be taken verbatim from the case where  $u \in \{w_1, \dots, w_{\tilde{d}}\}$  and  $u$  with the following modification. Remove every appearance of  $(a^\iota)^{\#u}$  and  $(\beta(u))^{\#u}$ . This concludes the proof of Claim 7.5.

It remain to show that

$$Q_{(w_1, \dots, w_{\tilde{d}})}(D) = \sum_{\beta \in \tilde{\mathcal{E}}^{[0]}} \prod_{p=1}^{\tilde{d}} \beta(w_p).$$

By definition it follows:

$$\begin{aligned} Q_{(w_1, \dots, w_{\tilde{d}})}(D) &= \llbracket \text{prod}(e_1, \dots, e_s) \rrbracket^{(D, \emptyset)} \\ &= \prod_{\substack{j \in [s] \\ y_j \notin V(\tilde{T}_{\overline{w}})}} \llbracket c_{z_j} \rrbracket^{(D, \emptyset)} \cdot \prod_{\substack{j \in [s] \\ z_j \in V(\tilde{T}_{\overline{w}})}} \llbracket \text{reg}_{\overline{w}}(z_j) \rrbracket^{(D, \emptyset)} \\ &\stackrel{\text{Claim 7.5}}{=} \prod_{\substack{j \in [s] \\ y_j \notin V(\tilde{T}_{\overline{w}})}} \left| \bigcup_{\iota_j \in \mathcal{L}_{z_j}^{[0]}} \tilde{\mathcal{E}}^{\iota_j} \right| \cdot \prod_{\substack{j \in [s] \\ z_j \in V(\tilde{T}_{\overline{w}})}} \sum_{\iota_j \in \mathcal{L}_{y_j}^{[0]}} \sum_{\beta \in \tilde{\mathcal{E}}^{\iota_j}} \prod_{q \in W_{z_j}} \beta(q)^{\#q} \end{aligned} \quad (7.3)$$

Furthermore, we have

$$\begin{aligned} \sum_{\beta \in \tilde{\mathcal{E}}^{[0]}} \prod_{p=1}^{\tilde{d}} \beta(w_p) &\stackrel{(*)}{=} \sum_{\beta \in \tilde{\mathcal{E}}^{[0]}} \prod_{q \in W_{v_{\text{root}}}} \beta(q)^{\#q} \\ &\stackrel{(**)}{=} \sum_{\iota_1 \in \mathcal{L}_{y_1}^{[0]}} \sum_{\beta_1 \in \tilde{\mathcal{E}}^{\iota_1}} \dots \sum_{\iota_s \in \mathcal{L}_{z_s}^{[0]}} \sum_{\beta_s \in \tilde{\mathcal{E}}^{\iota_s}} \prod_{q \in W_{v_{\text{root}}}} \beta(q)^{\#q} \\ &\quad \text{where } \beta := \beta_1 \cup \dots \cup \beta_s \\ &\stackrel{(***)}{=} \prod_{\substack{j \in [s] \\ z_j \notin V(\tilde{T}_{\overline{w}})}} \sum_{\iota_j \in \mathcal{L}_{z_j}^{[0]}} \sum_{\beta_j \in \tilde{\mathcal{E}}^{\iota_j}} \prod_{q \in W_{z_j}} \beta_j(q)^{\#q} \cdot \prod_{\substack{j \in [s] \\ z_j \in V(\tilde{T}_{\overline{w}})}} \sum_{\iota_j \in \mathcal{L}_{z_j}^{[0]}} \sum_{\beta_j \in \tilde{\mathcal{E}}^{\iota_j}} \prod_{q \in W_{w_j}} \beta_j(q)^{\#q} \\ &\stackrel{(***)}{=} \prod_{\substack{j \in [s] \\ z_j \notin V(\tilde{T}_{\overline{w}})}} \left| \bigcup_{\iota_j \in \mathcal{L}_{z_j}^{[0]}} \tilde{\mathcal{E}}^{\iota_j} \right| \cdot \prod_{\substack{j \in [s] \\ z_j \in V(\tilde{T}_{\overline{w}})}} \sum_{\iota_j \in \mathcal{L}_{z_j}^{[0]}} \sum_{\beta_j \in \tilde{\mathcal{E}}^{\iota_j}} \prod_{q \in W_{z_j}} \beta_j(q)^{\#q} \\ &\stackrel{(7.3)}{=} Q_{(w_1, \dots, w_{\tilde{d}})}(D) \end{aligned}$$

For Equation  $(*)$  note that every element  $z \in W_{v_{\text{root}}}$  occurs in the tuple  $\overline{w}$  exactly  $\#z$  times. In  $(**)$  we applied the decomposition lemma (Lemma 4.12). Every assignment  $\beta$  belongs to  $\tilde{\mathcal{E}}^{[0]}$  if and only if it can be decomposed into assignment  $\beta = \beta_1 \cup \dots \cup \beta_s$  such that there is an item  $\iota_j \in \mathcal{L}_{y_j}^{[0]}$  with  $\beta_j \in \tilde{\mathcal{E}}^{\iota_j}$ . In  $(***)$  we use the distributive law and the fact that  $\beta_j \subseteq \beta$  and  $\text{dom}(\beta_j) \supseteq W_{y_j}$ . In  $(****)$  we apply the fact from

### 7. An application concerning learning polynomial regression models over $q$ -hierarchical queries

Equation 7.2. Note that the condition for Equation 7.2 are given above. This concludes the correctness proof of  $Q_{(w_1, \dots, w_{\tilde{d}})}(D)$ .

We show now how to use these queries to obtain the result for Lemma 7.3. Let  $S := \{\bar{w} \in \text{vars}(Q)^{\tilde{d}}, \tilde{d} \leq 2d\}$ . Note that  $|S| \leq (|\text{vars}(Q)| + 1)^{2d}$ . We use in parallel for every query  $Q_{\bar{w}}$  for every  $\bar{w} \in S$  on the database  $D$  the data structure from Theorem 3.7. Note that  $t_a = O(1)$  for each such query.

For the initialisation routine we initialise the data structure for every query  $Q_{\bar{w}}$  for every  $\bar{w} \in S$ . This takes time  $O((|\text{vars}(Q)| + 1)^{2d}) = \text{poly}(Q)^{2d+1}$ .

Every time the database receives an update, we update each data structure for  $Q_{\bar{w}}$  for all  $\bar{w} \in S$ . This takes time  $O((|\text{vars}(Q)| + 1)^{2d} \text{poly}(Q)) = \text{poly}(Q)^{2d+1}$ .

When we receive as input a tuple  $\bar{w} \in S$  then we enumerate the tuple in time  $O(1)$ . Clearly,

$$Q_{(w_1, \dots, w_{\tilde{d}})}(D) = \sum_{\beta \in \tilde{\mathcal{E}}^{\{\emptyset\}}} \prod_{p=1}^{\tilde{d}} \beta(w_p)^{\#w_p}.$$

It takes time  $\text{poly}(Q)$  to output the value. This concludes the proof of Lemma 7.3.  $\square$

The proof of Theorem 3.8 follows from Lemma 7.3.

*Proof of Theorem 3.8.* Recall that for all  $(s_1, \dots, s_d) \in \mathcal{I}_d$  is

$$J(\xi)_{1, (s_1, \dots, s_d)} = \sum_{(j_1, \dots, j_d) \in \mathcal{I}_d} \theta_{(j_1, \dots, j_d)} \sum_{\beta \in \tilde{\mathcal{E}}^{\{\emptyset\}}} \prod_{p=1}^d \beta(x_{j_p}) \cdot \beta(x_{s_p}) - \sum_{\beta \in \tilde{\mathcal{E}}^{\{\emptyset\}}} \beta(y) \prod_{p=1}^d \beta(x_{j_p})$$

We use additionally the data structure of Lemma 7.3 for  $Q$  and  $D$  that can be initialised in time  $t_i = O((|\text{vars}(Q)| + 1)^{2d+1})$  and updated in time  $t_u = O(\text{poly}(Q)^{2d+1})$ . To receive the  $J(\xi)_{1, (s_1, \dots, s_d)}$  very fast, we compute  $\sum_{\beta \in \tilde{\mathcal{E}}^{\{\emptyset\}}} \beta(y) \prod_{p=1}^d \beta(x_{j_p})$  and  $\sum_{\beta \in \tilde{\mathcal{E}}^{\{\emptyset\}}} \prod_{p=1}^d \beta(x_{j_p}) \beta(x_{s_p})$ . Let  $\bar{w}^s$  ( $\bar{w}^j$ ) be the tuple, we obtain, if we remove every 0 from  $(s_1, \dots, s_d)$  ( $(j_1, \dots, j_d)$ ). Clearly,  $\sum_{\beta \in \tilde{\mathcal{E}}^{\{\emptyset\}}} \beta(y) \prod_{p=1}^d \beta(x_{j_p}) = Q_{\bar{w}^s}$  and  $\sum_{\beta \in \tilde{\mathcal{E}}^{\{\emptyset\}}} \prod_{p=1}^d \beta(x_{j_p}) \beta(x_{s_p}) = Q_{\bar{w}^j}$ . This can be done in time  $O(1)$  and to compute the complete jacobi-matrix it takes time  $|\mathcal{I}_d|^2 = \text{poly}(Q)^{2d+1}$ . Thus, it takes time  $\text{poly}(Q)^{2d+1}$  to do the forward pass.  $\square$

## 8. Outputting the $j$ th solution for $q$ -hierarchical conjunctive queries under updates

In this chapter, we show how the data structure for  $q$ -hierarchical conjunctive queries can be used to enhance for a given number  $j \in \mathbb{N}$ , the  $j$ th tuple in the enumeration with respect to a set of linear order  $\{\leq^v\}_{v \in V}$ . Here, we will achieve logarithmic (rather the constant) update time and answering time.

We will prove Theorem 3.9 that is restated in the following.

**Theorem 3.9.** *There is a dynamic algorithm that receives a  $q$ -hierarchical  $k$ -ary CQ with aggregates  $Q$  with aggregation time  $t_a$  and a  $\sigma$ -db  $D_0$  of size  $\|D_0\|$ , and computes within preprocessing time  $\text{poly}(Q) \cdot O(\|D_0\| \log(\|D_0\|))t_a$  a data structure that can be updated in time  $\text{poly}(Q) \cdot O(\log(\|D\|))t_a$  and allows the following:*

- (a) *output the query result size  $|Q(D)|$  in time  $O(1)$  and*
- (b) *enumerate the tuples in  $Q(D)$  in lexicographical order with delay  $\text{poly}(Q)$  and*
- (c) *test for an input tuple  $\bar{a} \in \mathbf{dom}^{k+\ell}$  if  $\bar{a} \in Q(D)$  within time  $t_t = \text{poly}(Q)$  and*
- (d) *upon input of a tuple  $\bar{b} \in \mathbf{dom}^{k+\ell}$  output the tuple*

$$\max \{ \bar{a} \in Q(D) : \bar{a} \leq \bar{b} \}$$

*if it exists, or a **SmallerThanMinimum**-message otherwise in time  $t_l = \text{poly}(Q)$  and*

- (e) *upon input of an arbitrary number  $j$ , take time  $\text{poly}(Q) \cdot O(\log(\|D\|))$  to immediately output the  $j$ th tuple that the enumeration procedure of (b) would output and*
- (f) *upon input of an arbitrary tuple  $\bar{a} \in Q(D)$ , take time  $\text{poly}(Q) \cdot O(\log(\|D\|))$  to immediately output the number  $j$  such that  $\bar{a}$  is the  $j$ th tuple the enumeration procedure of (b) would output,*

*where  $D$  is the current database.*

In the next lemma, we consider the special case for the  $j$ th problem and for the  $j$ th-reverse problem and for queries without aggregates. Using this lemma in combination with Lemma 4.31, we will receive that result in Theorem 8.1 by applying the reduction lemma (Lemma 6.18).

## 8. Outputting the $j$ th solution for $q$ -hierarchical conjunctive queries under updates

**Lemma 8.1.** *There is a dynamic algorithm that receives a  $q$ -hierarchical  $k$ -ary CQ  $Q$  (without aggregates) and a  $\sigma$ -db  $D_0$  of size  $\|D_0\|$  and computes within preprocessing time  $\text{poly}(Q) \cdot O(\|D_0\| \log(\|D_0\|))$  a data structure that can be updated in time  $\text{poly}(Q) \cdot O(\log(\|D\|))$  and allows the following:*

- (a) *output the query result size  $|Q(D)|$  in time  $O(1)$ ,*
- (b) *enumerate the tuples in  $Q(D)$  in lexicographical order with delay  $\text{poly}(Q)$ ,*
- (c) *test for an input tuple  $\bar{a} \in \text{dom}^{k+\ell}$  if  $\bar{a} \in Q(D)$  within time  $t_t = \text{poly}(Q)$ ,*
- (d) *upon input of an arbitrary number  $j$ , take time  $\text{poly}(Q) \cdot O(\log(\|D\|))$  to immediately output the  $j$ th tuple that the enumeration procedure of (b) would output if  $j \leq |Q(D)|$  and otherwise an **OutOfRange**-message,*
- (e) *upon input of an arbitrary tuple  $\bar{a} \in Q(D)$ , take time  $\text{poly}(Q) \cdot O(\log(\|D\|))$  to immediately output the number  $j$  such that  $\bar{a}$  is the  $j$ th tuple the enumeration procedure of (b) would output,*

where  $D$  is the current database.

A description of the data structure of Lemma 8.1 and a proof for part (d) is given in Section 8.1 and a proof for part (e) is given in Section 8.3.

In the Section 8.4 we give a proof for Theorem 3.9.

In the last section of this chapter (Section 8.5) we prove the lower bound stated in Theorem 3.16.

### 8.1. Output the $j$ th solution

This section is devoted to the proof of Lemma 8.1(d).

We start with an example. We consider the  $q$ -hierarchical query

$$Q = \{(x, y_1, z_1, y_2, z_2, z_3, y_3) : R(x, y_1, z_1) \wedge S(x, y_2, z_2) \wedge T(x, y_2, z_3) \wedge U(x, y_3)\}.$$

Let  $D$  be the database with

$$\begin{aligned} R^D &:= \{(1, 4, 1), (1, 4, 2), (1, 5, 1), (1, 5, 2), (2, 4, 9), (2, 4, 5), \\ &\quad (3, 1, 4), (3, 2, 5), (3, 3, 6)\} \\ S^D &:= \{(1, 1, 3), (1, 2, 5), (2, 1, 2), (3, 7, 7), (3, 7, 8)\} \\ T^D &:= \{(1, 1, 4), (1, 2, 6), (2, 1, 3), (3, 7, 9), (3, 7, 10)\} \\ U^D &:= \{(1, 7), (2, 8), (2, 9), (3, 1)\} \end{aligned}$$

Table 8.1 shows the tuples that will be output by the enumeration algorithm. See Figure 8.1 for a  $q$ -tree of  $Q$  and Figure 8.2 for an illustration of the data structure generated for  $Q$  and  $D$ . In Figure 8.2, the red number above an item  $i$  is the number  $|\tilde{\mathcal{E}}^i|$  (in Section 6.5 it is described how to compute such values).

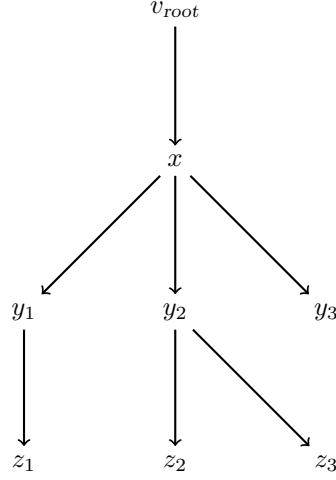
### 8.1. Output the $j$ th solution

Table 8.1.: Output of the enumeration of  $Q(D)$ .

No.	$x$	$y_1$	$z_1$	$y_2$	$z_2$	$z_3$	$y_3$
1	1	4	1	1	3	4	7
2	1	4	1	2	5	6	7
3	1	4	2	1	3	4	7
4	1	4	2	2	5	6	7
5	1	5	1	1	3	4	7
6	1	5	1	2	5	6	7
7	1	5	2	1	3	4	7
8	1	5	2	2	5	6	7
9	2	4	9	1	2	3	8
10	2	4	9	1	2	3	9
11	2	4	5	1	2	3	8
12	2	4	5	1	2	3	9
13	3	1	4	7	7	9	1
14	3	1	4	7	7	10	1
15	3	1	4	7	8	9	1
16	3	1	4	7	8	10	1
17	3	2	5	7	7	9	1
18	3	2	5	7	7	10	1
19	3	2	5	7	8	9	1
20	3	2	5	7	8	10	1
21	3	3	6	7	7	9	1
22	3	3	6	7	7	10	1
23	3	3	6	7	8	9	1
24	3	3	6	7	8	10	1

8. Outputting the  $j$ th solution for  $q$ -hierarchical conjunctive queries under updates

Figure 8.1.:  $q$ -tree for the example query

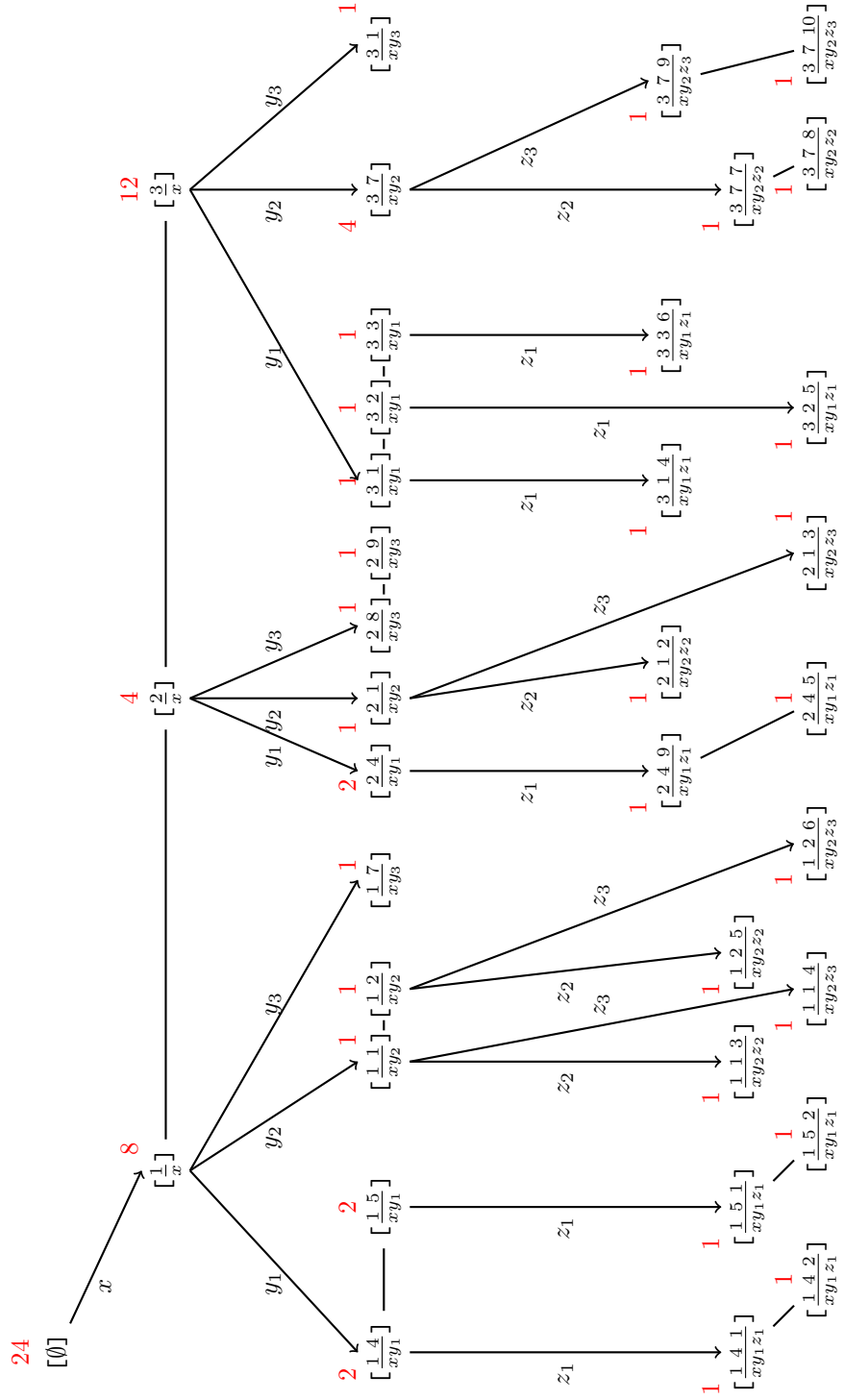


In this running example, we are now interested in the question: What is the  $j$ th tuple that will be output by the enumeration algorithm? Upon input of a number  $j$ , we want to quickly provide an answer to the question. For the rest of this example, let us assume that  $j = 20$ .

The value  $|\tilde{\mathcal{E}}^i|$  of the first item in the  $x$ -list of the *start*-item is 8. Therefore, we know that the enumeration procedure starts with 8 valuations  $\alpha$  where  $\alpha(x) = 1$ . Afterwards, the enumeration algorithm jumps to the next item in the start list to enumerate all the tuples where  $x$  is assigned to 2. As we know from its count value, there are 4 such tuples. Therefore, the component for  $x$  in the first 20 tuples in the enumeration are not assigned with 1 or 2. But in the enumeration there will follow 12 tuples where  $x$  is assigned to 3 and since  $8 + 4 < 20 \leq 8 + 4 + 12$  (these values are taken from the  $|\tilde{\mathcal{E}}^i|$  values from the items in the  $x$ -list of  $[\emptyset]$ ), we know that  $x$  must be assigned with 3. Note that the 20th tuple in the enumeration is exactly the  $20 - 12 = 8$ th tuple in the enumeration from  $[\frac{3}{x}]$ . So, we try to find the 8th tuple in the enumeration from  $[\frac{3}{x}]$ .

Now, we continue with the item  $[\frac{3}{x}]$  and check in its lists how many assignments they will enumerate. We start with the  $y_3$ -list (on the right) and we see that there is only one assignment, namely  $\beta(y_3) = 1$ . Now, we know that for all 12 tuples which are enumerated from  $[\frac{3}{x}]$ , the variable  $y_3$  will be assigned with 1. Then, we consider the  $y_2$ -list in the middle and we see that it has 4 assignments to enumerate. The assignments are  $\beta_1 = \frac{7 \ 7 \ 9}{y_2 \ z_2 \ z_3}, \beta_2 = \frac{7 \ 7 \ 10}{y_2 \ z_2 \ z_3}, \beta_3 = \frac{7 \ 8 \ 9}{y_2 \ z_2 \ z_3}, \beta_4 = \frac{7 \ 8 \ 10}{y_2 \ z_2 \ z_3}$ . Now, in the first tuple of the enumeration from  $[\frac{3}{x}]$  the variables  $y_2, z_2, z_3$  are assigned with the values of  $\beta_1$ , the second tuple with the values of  $\beta_2$ , the third tuple with the values of  $\beta_3$ , the fourth tuple with the values of  $\beta_4$  and this repeats for three times. To find out

Figure 8.2.: Illustration of the data structure of the example database.



8.1. Output the  $j$ th solution

## 8. Outputting the $j$ th solution for $q$ -hierarchical conjunctive queries under updates

which is the right assignment for tuple number 8, we simply compute 8 modulo 4. We receive 0 and then we know that  $\beta_4$  is the correct assignment, since for every 4th tuple in the enumeration the variables  $y_3, z_2, z_3$  are assigned with the values of  $\beta_4$ . Finally, we have to take into account the  $y_1$ -list. It has the three assignments  $\gamma_1 = \frac{1}{y_1 z_1}, \gamma_2 = \frac{2}{y_1 z_1}, \gamma_3 = \frac{3}{y_1 z_1}$ . The enumeration will output for 4 times  $\gamma_1$  (in combination with every  $\beta_j$ ), then output for 4 times  $\gamma_2$  and afterwards 4 times  $\gamma_3$ . It is straightforward to see that  $\gamma_2$  is the assignment, we are looking for. To compute the right solution we divide 8 by 4 and we receive the information that we need to output  $\gamma_2$ . (Note that we have to round up the solution of the division). Thus, we obtain the solution tuple  $(3, 2, 5, 7, 8, 10, 1)$ .

By using the idea illustrated in this example, we can enrich the data structure in Theorem 3.3 to obtain a proof for Lemma 8.1.

For the remainder of this section we assume that  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$  is a  $q$ -hierarchical conjunctive query,  $\text{vars}(Q) = \{x_1, \dots, x_m\}$  with  $0 \leq k \leq m$ , and  $Q$  is of the form

$$Q = \{(x_1 \dots x_k, b_1, \dots, b_\ell) : \exists x_{k+1} \dots \exists x_m (\psi_1 \wedge \dots \wedge \psi_d)\}, \quad (8.1)$$

where  $b_1, \dots, b_\ell \in \mathbf{dom}$  and  $\psi_1, \dots, \psi_d$  are atomic queries of schema  $\sigma$ .

We need that data structure from Section 6.6 for the counting query

$$Q_c := \{(\text{prod}(c_{x_1}, \dots, c_{x_m})) : \varphi\}$$

from Lemma 6.11 available.

We need the following variants of the modulo operation and the division operation. For all integers  $a, b \in \mathbb{N}$  let

$$a \widehat{\text{mod}} b := \begin{cases} b & \text{if } a \bmod b \equiv 0 \\ a \bmod b & \text{otherwise} \end{cases}$$

and

$$a \widehat{\text{div}} b := \lceil a/b \rceil.$$

To compute for a given number  $j \in \mathbb{N}$  the  $j$ th tuple of the enumeration of  $Q(D)$ , we start the procedure GETASSIGN of Algorithm 11 with  $j$  and the *start*-item as arguments.

To implement line 10 we use the following lemma:

**Lemma 8.2.** *There is a dynamic algorithm that receives as input a list  $\langle n_1, \dots, n_s \rangle$  of type  $\mathbb{N}$  and computes within  $t_p = O(s \cdot \log s)$  preprocessing time a data structure that can be updated in time  $t_u = O(\log s)$  and allows on input of a number  $m \in \mathbb{N}$  either*

- *output the index  $p \in [s]$  such that  $\sum_{i=1}^{p-1} n_i < m \leq \sum_{i=1}^p n_i$  and the number  $\sum_{i=1}^{p-1} n_i$  if  $m \leq \sum_{i=1}^s n_i$  or*
- *output the message *OutOfRange* if  $m > \sum_{i=1}^s n_i$ .*



---

**Algorithm 11** Algorithm to output the  $j$  tuple.

---

```

1: function GETASSIGN( $j, i, \{<^v\}_{v \in \text{succ}(v^i)}$ )
2:   Input: Number  $j \in [\left|\tilde{\mathcal{E}}^i\right|]$ , item  $i$ 
3:   Output: The assignment for the  $j$ th tuple of enumerating  $\tilde{\mathcal{E}}^i$ .
4:   if  $v^i$  is a leaf in the  $q$ -tree then return  $\alpha^i$ 
5:   Let  $u_1 <^{v^i} \dots <^{v^i} u_d$  be the children of  $v^i$  in  $T_{\text{free}}$ .
6:   for all  $p \in [d]$  do
7:      $c_p \leftarrow \prod_{q=p+1}^d \sum_{\iota \in \mathcal{L}_{u_q}^i} \left|\tilde{\mathcal{E}}^\iota\right|$ 
8:      $m_p = j \widehat{\text{div}} c_p$ 
9:      $n_p = m_p \widehat{\text{mod}} \sum_{\iota \in \mathcal{L}_{u_q}^i} \left|\tilde{\mathcal{E}}^\iota\right|$ 
10:    For the list  $\langle i_1, \dots, i_s \rangle = \mathcal{L}_u^i$  compute the  $q \in [d]$  such that  $\sum_{r=1}^{q-1} \left|\tilde{\mathcal{E}}^{i_r}\right| <$ 
        $n_p \leq \sum_{r=1}^q \left|\tilde{\mathcal{E}}^{i_r}\right|$ .
11:    Let  $\hat{j}$  be  $n_p - \sum_{r=1}^{q-1} \left|\tilde{\mathcal{E}}^{i_r}\right|$ .
12:     $\alpha_p \leftarrow \text{GETASSIGN}(\hat{j}, i_q, \{<^v\}_{v \in \text{succ}(u_p)})$ 
13:   return  $\alpha_1 \cup \dots \cup \alpha_d$ 

```

---

and upon input of an item  $n_r$  from the list, output the number  $\sum_{i=1}^r n_i$ .

Note an update means here, an update on a list as described in page 13.

A proof for this lemma is given in Section 8.2.

To realise the algorithm, we need to get fast access to the values  $\left|\tilde{\mathcal{E}}^i\right|$  and to  $\sum_{\iota \in \mathcal{L}_u^i} \left|\tilde{\mathcal{E}}^\iota\right|$  for every child  $u$  of  $v^i$  in  $T_{\text{free}}$ . This can be done by using the data structure we obtain from the aggregate query that outputs the number of results in Section 6.5. Note that the data structure stores the numbers for each item on the item. Furthermore, we store for every item  $i$  and for every child  $u$  of  $v^i$  a list that contains the values  $\langle \left|\tilde{\mathcal{E}}^\iota\right| : \iota \in \mathcal{L}_u^i \rangle$  and the data structure from Lemma 8.2. The order of the elements in the list corresponds to the order of the items in the list, i.e., the  $m$ th element in the list is  $\left|\tilde{\mathcal{E}}^\iota\right|$  if and only if  $\iota$  is the  $m$ th element in  $\mathcal{L}_u^i$ . The data structure can be used to get fast access to  $q$  in line 10 and to compute the value  $\hat{j}$  in line 11.

If we update the database we insert/delete a constant number of items and modify a constant number of  $\left|\tilde{\mathcal{E}}^i\right|$  values. For every item we modify in the data structure, we update the corresponding element in the data structure with Lemma 8.2. In particular, we need  $\text{poly}(Q) \times O(\log(\|D\|))$  since we have to update  $\text{poly}(Q)$  lists in  $O(\log(\|D\|))$  time.

We show now that Algorithm 8.1 is correct.

**Claim 8.3.** Upon input of an item  $i$  and  $j \in [\left|\tilde{\mathcal{E}}^i\right|]$  and a family of linear orderings

## 8. Outputting the $j$ th solution for $q$ -hierarchical conjunctive queries under updates

$\{<^v\}_{v \in \text{succ}(v^i)}$  the function  $\text{GETASSIGN}(j, i, \{<^v\}_{v \in \text{succ}(v^i)})$  in Algorithm 11 outputs the  $j$ th assignment yielded by the function  $\text{ENUM}(i, \{<^v\}_{v \in \text{succ}(v^i)})$  in Algorithm 3. Furthermore, Algorithm 11 takes time  $O(|\text{succ}(v^i)| (|\text{succ}(v^i)| + \log(|\text{adom}(D)|)))$ .

*Proof.* We show this claim by induction over the height of an item in  $T_{\text{free}}$ .

For the induction base, let us consider an item of height 0. Then,  $v^i$  is a leaf in  $T_{\text{free}}$  and  $\tilde{\mathcal{E}}^i = \{\alpha^i\}$  and  $j = 1 = |\tilde{\mathcal{E}}^i|$  and the function in  $\text{GETASSIGN}$  outputs  $\alpha^i$ , that is the unique assignment yielded by  $\text{ENUM}(i, \emptyset)$ . This takes time  $O(1)$ .

For the induction step let us consider an item of height  $h$  and let  $u_1 <^{v^i} \dots <^{v^i} u_d$  be the children of  $v^i$  in  $T_{\text{free}}$  and  $<^{v^i}$  be a linear order we receive as input. Note by the decomposition lemma (Lemma 4.12) the  $j$ th assignment can be decomposed to  $\alpha_1 \cup \dots \cup \alpha_d$  where  $\alpha_p \in \tilde{\mathcal{E}}^{\iota_p}$  with  $\iota_p \in \mathcal{L}_{u_p}^i$  for all  $p \in [d]$ . For all  $p \in [d]$  it holds that  $\alpha_p$  is taken as an element from the loops **for**  $\iota_p \in \mathcal{L}_{u_p}^i$  **do for**  $\alpha_p \in \text{ENUM}(\iota_p, \{<^v\}_{v \in \text{succ}(v^{\iota_p})})$  **do**. By construction of the algorithm, the number of assignments  $\beta$  with  $\beta \supseteq \alpha_p$  that will be yielded is  $c_p$ . Thus, the assignment we need is in the  $m_p$ th iteration. Note that there might be for-loops that force to restart the for-loops. But, the for-loops mentioned above have exactly  $\sum_{\iota \in \mathcal{L}_{u_q}^i} |\tilde{\mathcal{E}}^\iota|$  iterations. Thus, the assignment we need, is the  $n_p$ th assignment in the iteration of the for-loops mentioned above. In the next lines, the algorithm picks from the for-loops the item  $i_q$  that will be supposed in the  $n_p$ th iteration. Then, we subtract the number of iterations with the item that appears before  $i_q$  and take this number to receive the assignment  $\alpha_p$  from  $\text{GETASSIGN}(\hat{j}, i_q, \{<^v\}_{v \in \text{succ}(u_p)})$  that works by induction hypothesis correctly. Note that  $1 \leq n_p \leq |\tilde{\mathcal{E}}^{\iota_p}|$ .

Line 7 takes time  $O(d)$  and line 10 and 11 takes  $O(\log(|\text{adom}(D)|))$  time, since we use Lemma 8.2 to compute the values. By induction hypothesis line 12 takes time  $O(|\text{succ}(u_p)| (|\text{succ}(u_p)| + \log(n)))$ . The other lines in the body of the for-loop takes time  $O(1)$ . Since we use the data structure from Section 6.6 for the query  $Q_c$  we have constant access to the  $|\tilde{\mathcal{E}}^i|$  values. All in all, the time needed for the whole for-loop is

$$\begin{aligned} & \sum_{p \in [d]} O(d) + O(\log(|\text{adom}(D)|)) + O(|\text{succ}(u_p)| (|\text{succ}(u_p)| + \log(|\text{adom}(D)|))) \\ &= O(|\text{succ}(v^i)| (|\text{succ}(v^i)| + \log(|\text{adom}(D)|))) \end{aligned}$$

Therefore, it takes time  $O(|\text{succ}(v^i)| (|\text{succ}(v^i)| + \log(|\text{adom}(D)|)))$  to return the assignment.  $\square$

To conclude the proof of Lemma 8.1 we check upon input of a number  $j \in \mathbb{N}$  if  $j \leq |\varphi(D)|$  using the data structure for output the number  $|\varphi(D)|$ . This can be done in  $O(1)$ . If this is not the case, we output the **OutOfRange** message. Otherwise, we use Algorithm 11 with input  $j$  and the start item  $[\emptyset]$  to obtain the  $j$ th assignment  $\alpha$  of enumeration of  $\varphi(D)$  in time  $O(|\text{succ}(v^i)| (|\text{succ}(v^i)| + \log(n)))$  and output the tuple  $(\alpha(x_1), \dots, \alpha(x_k), b_1, \dots, b_\ell)$ . This concludes the proof of Lemma 8.1(d).

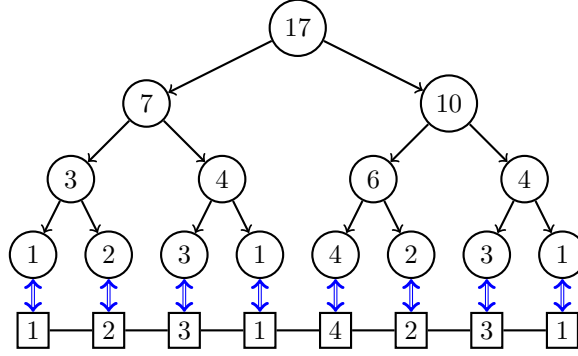


Figure 8.3.: Example of a sumtree.

## 8.2. Proof of Lemma 8.2

Let  $\mathcal{L}$  be a list of integers of size  $s$ . The data structure we use is an enrichment of the data structure in Theorem 2.1, i.e., for the list  $\mathcal{L}$  we have an balanced tree available, where every leaf is associated with an item in the list and the leafs are sorted as the elements in the list. Additionally, we store for every vertex  $v$  a field  $\text{num}(v)$  that is defined as follows: For all leafs  $v$  the number  $\text{num}(v)$  is the number of the item in the list it is indicated to. For all internal vertices  $v$  the number  $\text{num}(v)$  is the sum of the number values of its children. For all vertices  $v$ , we let  $\mathcal{F}_v \subseteq [s]$  be the set of indices  $j$  such that there is a path from  $v$  to a leaf that indicates  $j$ . Note that it easily follows that  $\text{num}(v) = \sum_{j \in \mathcal{F}_v} \mathcal{L}[j]$ . With  $\mathcal{L}[j]$  we denote the  $j$ th element in the list  $\mathcal{L}$ . See Figure 8.3 for the balanced tree for a list, where the numbers in the nodes  $v$  depict the number  $\text{num}(v)$ .

We show on an example: upon input of the number 12 and the data structure in Figure 8.3, we show how to find the index  $p$  such that  $\sum_{i=1}^{p-1} n_i < 12 \leq \sum_{i=1}^p n_i$ . First of all we check on the root  $v$  whether its number is greater than 12. Since it is not the case, we now consider the root and check the number of its children. As we see, the left child  $v_{\text{left}}$  has the number 7. Since  $7 = \sum_{j \in \mathcal{F}_{v_{\text{left}}}} \mathcal{L}[j]$ , it follows that the sum of the first 4 elements in the list is 7. Thus, the index  $p$ , we search, must be greater than 4. In such a case we continue with the right vertex  $v_{\text{right}}$  and its left child  $w_{\text{left}}$  and its right child  $w_{\text{right}}$ . We also store the number 7 from  $v_{\text{left}}$  in a variable **current**. We consider now  $w_{\text{left}}$  and see that  $\text{current} + \text{num}(w_{\text{left}}) = 7 + 6 = 13 > 12$ . In such a case we know that the sum of the first six elements in the list is greater than 12 and therefore for the  $p$  we search is  $p \leq 6$ . Note that  $\text{num}(w_{\text{left}})$  is the sum  $\mathcal{L}[5] + \mathcal{L}[6]$ . Then, we jump in the next step to the node  $w_{\text{left}}$  and consider its left child  $u_{\text{left}}$  and right child  $u_{\text{right}}$ . We do not change the **current** value. It still stores the sum of the values of the first four elements in the list. Again, we add **current** and  $\text{num}(u_{\text{left}})$  to

## 8. Outputting the $j$ th solution for $q$ -hierarchical conjunctive queries under updates

get the sum of the first five elements in the list 11. Since

$$\text{current} + \text{num}(w_{\text{left}}) = \sum_{i=1}^5 \mathcal{L}[i] < 12 \leq \sum_{i=1}^6 \mathcal{L}[i] = 13 = \text{current} + \text{num}(w_{\text{left}})$$

we have found the searched index  $p = 6$ .

Algorithm 8.2 shows how to solve this problem with the data structure. First of all, we check if the input number  $n$  is greater than the number value of the root. Note that the number on the root is  $\sum_{i=1}^s \mathcal{L}[i]$ . In that case, we output the **OutOfRange** message. Otherwise we start the procedure **FIND** in Algorithm 12 with input values 0 for **current**,  $m$  for  $j$  and the root node for  $v$ .

---

**Algorithm 12** Procedure to find the index.

---

```

1: function FIND(current,  $m$ ,  $v$ )
2:   Input: current,  $m \in \mathbb{N}$  and vertex  $v$ 
3:   if  $v$  is a leaf then
4:     return  $p$  is the index of element indicated by  $v$  and current.
5:   Let  $v_{\text{left}}$  be the left child of  $v$ 
6:   Let  $v_{\text{right}}$  be the right child of  $v$ 
7:   if  $m \leq \text{current} + \text{num}(v_{\text{left}})$  then
8:     return FIND(current,  $m$ ,  $v_{\text{left}}$ )
9:   else
10:    return FIND(current +  $\text{num}(v_{\text{left}})$ ,  $m$ ,  $v_{\text{right}}$ )

```

---

For all vertices  $v$  let  $\text{first}(v) := \min \mathcal{F}_v$  and  $\text{last}(v) := \max \mathcal{F}_v$ . To show the correctness of the algorithm, we show by induction of the recursion depth of **FIND** that by calling the **FIND**(**current**,  $m$ ,  $v$ ) procedure the following invariant holds:

•

$$\sum_{i=1}^{\text{first}(v)-1} \mathcal{L}[i] < m \leq \sum_{i=1}^{\text{last}(v)} \mathcal{L}[i]$$

•

$$\text{current} = \sum_{i=1}^{\text{first}(v)-1} \mathcal{L}[i]$$

*Induction base  $d = 1$ .* When we call **FIND**(0,  $m$ ,  $v_{\text{root}}$ ) at the beginning, we have made sure that  $m \leq \sum_{i=1}^s \mathcal{L}[i]$ . Note that since  $\text{first}(v_{\text{root}}) = 1$  and  $\text{last}(v_{\text{root}}) = s$ , we have

$$\sum_{i=1}^{\text{first}(v_{\text{root}})-1} \mathcal{L}[i] = 0 < m \leq \sum_{i=1}^s \mathcal{L}[i] = \sum_{i=1}^{\text{last}(v_{\text{root}})} \mathcal{L}[i]$$

and

$$\text{current} = \sum_{i=1}^{\text{first}(v_{\text{root}})-1} \mathcal{L}[i] = 0 .$$

## 8.2. Proof of Lemma 8.2

*Inductive step  $d-1 \rightarrow d$ .* We call  $\text{FIND}(\text{current}, m, w)$  during the execution of  $\text{FIND}(\text{current}', m, u)$  that have called by the algorithm in recursion depth  $d-1$ . By induction hypothesis we have

$$\sum_{i=1}^{\text{first}(u)-1} \mathcal{L}[i] < m \leq \sum_{i=1}^{\text{last}(u)} \mathcal{L}[i]$$

and

$$\text{current}' = \sum_{i=1}^{\text{first}(u)-1} \mathcal{L}[i]$$

- **Case 1:**  $\text{FIND}(\text{current}, m, w)$  was called by  $\text{FIND}(\text{current}', m, u)$  since  $m \leq \text{current}' + \text{num}(w)$  where  $w$  is the left child of  $u$ . Then  $\text{current} = \text{current}'$ . Since  $w$  is the left child of  $u$ , it follows that  $\text{first}(w) = \text{first}(u)$  and therefore we have

$$\text{current} = \text{current}' = \sum_{i=1}^{\text{first}(u)-1} \mathcal{L}[i] = \sum_{i=1}^{\text{first}(w)-1} \mathcal{L}[i]$$

and

$$\begin{aligned} \sum_{i=1}^{\text{first}(w)-1} \mathcal{L}[i] &= \sum_{i=1}^{\text{first}(u)-1} \mathcal{L}[i] < m \leq \text{current}' + \text{num}(w) \\ &= \sum_{i=1}^{\text{first}(w)-1} \mathcal{L}[i] + \sum_{i=\text{first}(w)}^{\text{last}(w)} \mathcal{L}[i] = \sum_{i=1}^{\text{last}(w)} \mathcal{L}[i] \end{aligned}$$

- **Case 2:**  $\text{FIND}(\text{current}, m, w)$  was called by  $\text{FIND}(\text{current}', m, u)$  since  $m > \text{current}' + \text{num}(v)$  where  $v$  is the left child of  $u$ . Then  $\text{current} = \text{current}' + \text{num}(v)$  and  $w$  is the right child of  $u$ . Since  $w$  is the right child of  $u$ , it follows that  $\text{last}(w) = \text{last}(u)$ . Therefore, we have

$$m \leq \sum_{i=1}^{\text{last}(u)} \mathcal{L}[i] = \sum_{i=1}^{\text{last}(w)} \mathcal{L}[i]$$

By

$$\begin{aligned} m > \text{current}' + \text{num}(v) &= \sum_{i=1}^{\text{first}(u)-1} \mathcal{L}[i] + \sum_{i=\text{first}(v)}^{\text{last}(v)} \mathcal{L}[i] \\ \stackrel{\text{first}(v) \equiv \text{first}(u)}{\sum_{i=1}^{\text{last}(v)} \mathcal{L}[i]} &= \sum_{i=1}^{\text{first}(w)-1} \mathcal{L}[i] \end{aligned}$$

and

$$\text{current} = \text{current}' + \text{num}(v) = \sum_{i=1}^{\text{first}(w)-1} \mathcal{L}[i]$$

8. *Outputting the  $j$ th solution for  $q$ -hierarchical conjunctive queries under updates*

it follows that the invariant holds before  $\text{FIND}(\mathbf{current}, m, w)$  is called.

This concludes the induction.

Let  $v$  be a leaf and let  $p$  be the index of the element indicated by  $v$ . Then, we have  $\text{first}(v) = \text{last}(v) = p$ . In particular, we have

$$\sum_{i=1}^{\text{first}(v)-1} \mathcal{L}[i] = \sum_{i=1}^{p-1} \mathcal{L}[i] < m \leq \sum_{i=1}^{\text{last}(v)} \mathcal{L}[i] = m = \sum_{i=1}^p \mathcal{L}[i]$$

and

$$\mathbf{current} = \sum_{i=1}^{\text{first}(v)-1} \mathcal{L}[i] = \sum_{i=1}^{p-1} \mathcal{L}[i]$$

This is exactly the index and  $\mathbf{current}$  is the value we are searching for and this concludes the proof of the correctness of Algorithm 8.2.

To implement the routine that upon input of an item  $n_r$  from the list, outputs the number  $\sum_{i=1}^r n_i$ , we use Algorithm 13 on input of the leaf  $w$  in the sum tree that is associated to  $n_r$  and  $\mathbf{sum} = \mathbf{num}(w)$ .

---

**Algorithm 13** Procedure to compute the sum  $\sum_{i=1}^r n_i$ .

---

```

1: function GETSUM( $\mathbf{current}, m, w$ )
2:   Input:  $\mathbf{current} \in \mathbb{N}$  and vertex  $w$ .
3:   if  $w$  is the root of the sum tree then
4:     return  $\mathbf{current}$ 
5:   else
6:     Let  $u$  be the parent of  $w$ .
7:     if  $w$  is the left child of  $u$  or the only child of  $u$  then
8:       GETSUM( $\mathbf{current}, w$ ).
9:     else
10:      Let  $u_{\text{left}}$  be the left child of  $u$ .
11:      GETSUM( $\mathbf{current} + \mathbf{num}(u_{\text{left}}), v$ ).

```

---

For all vertices  $v$  let  $\text{first}(v) := \min \mathcal{F}_v$  and  $\text{last}(v) := \max \mathcal{F}_v$ . To show the correctness of the algorithm, we prove the following. When we start  $\text{GETSUM}(\mathbf{num}(v), v)$  where  $v$  is a leaf in the sum tree, then in every recursive step the following invariant holds. If  $\text{GETSUM}(\mathbf{current}, u)$  will be executed, we have

$$\mathbf{current} = \sum_{i=\text{first}(u)}^{\text{last}(v)} \mathcal{L}[i].$$

We show this by an induction over the recursive depth. For the induction base let us consider the case that we start with  $\text{GETSUM}(\mathbf{num}(v), v)$ . Then the claim

holds since  $v$  is a leaf and thus

$$\mathbf{current} = \mathbf{num}(v) = \sum_{i=\mathit{first}(v)}^{\mathit{last}(v)} \mathcal{L}[i].$$

For the inductive step, let us consider a recursive step with depth  $d + 1$ .

- *Case 1:*  $\mathbf{GETSUM}(\mathbf{current}, w)$  was executed in  $\mathbf{GETSUM}(\mathbf{current}, u)$  (with recursive depth  $d$ ) in line 8. Then  $w$  is the left child of  $u$  or the only child of  $u$  and thus  $\mathit{first}(u) = \mathit{first}(w)$ . We obtain from the induction hypothesis the following.

$$\mathbf{current} = \sum_{i=\mathit{first}(u)}^{\mathit{last}(v)} \mathcal{L}[i] = \sum_{i=\mathit{first}(w)}^{\mathit{last}(v)} \mathcal{L}[i].$$

- *Case 2:*  $\mathbf{GETSUM}(\mathbf{current}, w)$  was executed in  $\mathbf{GETSUM}(\mathbf{current}, u)$  (with recursive depth  $d$ ) in line 11. Then  $w$  is the right child of  $u$  and thus  $\mathit{first}(w) = \mathit{last}(u_{\text{left}}) + 1$  where  $u_{\text{left}}$  is the left child of  $u$ . We obtain from the induction hypothesis and the definition of  $\mathbf{num}(u_{\text{left}})$  the following.

$$\begin{aligned} \mathbf{current} + \mathbf{num}(u_{\text{left}}) &= \sum_{i=\mathit{first}(w)}^{\mathit{last}(v)} \mathcal{L}[i] + \sum_{i=\mathit{first}(u_{\text{left}})}^{\mathit{last}(u_{\text{left}})} \mathcal{L}[i] \\ &= \sum_{i=\mathit{last}(u_{\text{left}})+1}^{\mathit{last}(v)} \mathcal{L}[i] + \sum_{i=\mathit{first}(u_{\text{left}})}^{\mathit{last}(u_{\text{left}})} \mathcal{L}[i] \\ &= \sum_{i=\mathit{first}(u_{\text{left}})}^{\mathit{last}(v)} \mathcal{L}[i]. \end{aligned}$$

Thus, it follows that the number returned by  $\mathbf{GETSUM}(\mathbf{num}(v), v)$  where  $v$  is the leaf for  $n_r$  is equal to  $\sum_{i=1}^r n_i$ . Since the recursion depth is bounded by the depth of the tree which is  $O(\log(s))$ , Algorithm 13 takes time  $O(\log(s))$ .

We show now how to maintain sum trees under updates.

Recall, that the data structure is an enrichment of the data structure of Theorem 2.1. In particular, the tree with the sums is a AVL-tree and during an update operation, one has to rotate the tree to get a balanced tree (see Figures 2.1 and 2.2). We enrich the implementation of the update steps as follows.

- When we insert or delete a node to/from the list, we update the  $\mathbf{num}(v)$  value to all nodes on the path from the root to the corresponding leaf. This takes  $O(s)$  time.
- After a closer look to the rotations operations (see [65]), we see that an rotation changes the children of at most three nodes. We change the algorithm for updating such that we recompute the number  $\mathbf{num}(v)$  values. The operations for such an update takes time  $O(1)$ .

### 8. Outputting the $j$ th solution for $q$ -hierarchical conjunctive queries under updates

Therefore updating the data structure can be done in  $O(\log(s))$ .

This concludes the proof of Lemma 8.2.

## 8.3. $j$ th reverse

In this section, we consider the  $j$ th reverse problem, i.e., we prove Lemma 8.1(e). For the remainder of this section we assume that  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$  is a  $q$ -hierarchical conjunctive query,  $\text{vars}(Q) = \{x_1, \dots, x_m\}$  with  $0 \leq k \leq m$ , and  $Q$  is of the form

$$Q = \{(x_1 \dots x_k, b_1, \dots, b_\ell) : \exists x_{k+1} \dots \exists x_m (\psi_1 \wedge \dots \wedge \psi_d)\}, \quad (8.2)$$

where  $b_1, \dots, b_\ell \in \mathbf{dom}$  and  $\psi_1, \dots, \psi_d$  are atomic queries of schema  $\sigma$ . Let  $Q$  be a  $q$ -hierarchical query and let  $D$  be a database. We receive as input an assignment  $\alpha$  and we request the number  $j$  such that  $(\alpha(x_1), \dots, \alpha(x_k), b_1, \dots, b_\ell)$  is the  $j$ th tuple in the enumeration of  $Q(D)$ .

To compute the number, we define the number  $\text{jth}(v)$  for all  $v \in \text{vars}(Q)$  such that  $\text{jth}(v)$  is the number of steps until we output the assignment  $\alpha|_{\text{vpath}[v]} \cup \alpha|_{\text{succ}(v)}$  in the procedure  $\text{Enum}([\alpha|_{\text{vpath}[v]}], \{<^w\}_{w \in \text{succ}(v|_{\text{vpath}[v]})})$ . The number of steps means here number of assignments that will be yielded. The number  $\text{jth}(v)$  can be defined inductively over the height of  $v$ .

- If  $v$  has height 0, we have  $\text{jth}(v) = 0$ .
- If  $v$  has height  $h > 0$ : Let  $u_1, \dots, u_s$  be the children of  $v$  in  $T_{\text{free}}$  and for all  $j \in [s]$  let  $\mathcal{I}_j$  be the set of items that appear before  $[\alpha|_{\text{vpath}[u_j]}]$  in the  $u_j$ -list of  $[\alpha|_{\text{vpath}[v]}]$ . We count the number of steps in the procedure  $\text{Enum}([\alpha|_{\text{vpath}[v]}], \{<^w\}_{w \in \text{succ}(v|_{\text{vpath}[v]})})$  until we receive  $\alpha_i = \alpha|_{\text{succ}(u_i)}$  for all  $i \in [s]$ . By construction of the algorithm, we run the first two for-loops until  $\alpha_1 = \alpha|_{\text{succ}(u_1)}$ , then the next two for-loops until  $\alpha_2 = \alpha|_{\text{succ}(u_2)}$  and so on, i.e., the number can be expressed as

$$\text{jth}(v) = \sum_{j=1}^s \text{“number of steps until } \alpha_j = \alpha|_{\text{succ}(u_j)} \text{”}.$$

To compute the number of steps until  $\alpha_j = \alpha|_{\text{succ}(u_j)}$  for all  $j \in [s]$ , we count the number of steps until we reach one step before  $\alpha_j = \alpha|_{\text{succ}(u_j)}$  for all  $j \in [s]$ . We add then 1 to obtain the number of steps until  $\alpha_j = \alpha|_{\text{succ}(u_j)}$ . We need  $\sum_{i \in \mathcal{I}_j} |\tilde{\mathcal{E}}^i|$  steps until the for loop **for**  $\hat{l}_j \in \mathcal{L}_{u_j}^i$  **do** reaches  $\hat{l}_j = [\alpha|_{\text{vpath}[u_j]}]$ . Then, by induction hypothesis, we need  $\text{jth}(u_j)$  steps until the for loop **for**  $\alpha_j \in \text{ENUM}(\hat{l}_j, \{<^v\}_{v \in \text{succ}(u_j)})$  **do** reaches  $\alpha_j = \alpha|_{\text{succ}(u_j)}$ . For every iteration we have



$\prod_{j=j+1}^s \sum_{\iota \in \mathcal{L}_{[\alpha]_{\text{vpath}[v]}^{u_j}}^{u_j}} |\tilde{\mathcal{E}}^\iota|$  steps that run over the other for loops. We obtain

$$\text{jth}(v) = \sum_{j=1}^s \left( \left( \sum_{\iota \in \mathcal{I}_j} |\tilde{\mathcal{E}}^\iota| + \text{jth}(u_j) - 1 \right) \cdot \prod_{\hat{j}=j+1}^s \sum_{\iota \in \mathcal{L}_{u_{\hat{j}}}^{[\alpha]_{\text{vpath}[v]}^{u_{\hat{j}}}}} |\tilde{\mathcal{E}}^\iota| + 1 \right)$$

To compute the numbers  $\text{jth}(v)$  we use a bottom up algorithm on the  $q$ -tree of  $Q$ . The numbers  $\sum_{\iota \in \mathcal{I}_j} |\tilde{\mathcal{E}}^\iota|$  and  $\sum_{\iota \in \mathcal{L}_{u_j}^{[\alpha]_{\text{vpath}[v]}^{u_j}}} |\tilde{\mathcal{E}}^\iota|$  can be computed from the sum tree described in the first section of this chapter. Note that we have for every item  $i$  and for every  $u \in \text{child}(v^i)$  a list with a sum tree available that contains for every  $\iota \in \mathcal{L}_u^i$  the number  $|\tilde{\mathcal{E}}^\iota|$ . From Lemma 8.2 it follows that there is a routine where we can output the sum  $\sum_{\iota \in \mathcal{I}_j} |\tilde{\mathcal{E}}^\iota|$  in time  $O(\log(\text{adom}(D)))$ .

All in all, we need  $O(\text{poly}(Q) \log(\text{adom}(D)))$  time to compute the  $\text{jth}(v)$  values for all  $v \in \text{vars}(Q)$ . This concludes the proof of Lemma 8.1(e).

## 8.4. Proof for Theorem 3.9

In the following corollary we combine data structure of Lemma 4.31 with the data structure of Lemma 8.1. Recall that the data structure of Lemma 4.31 is simply a lexicographically ordered data structure. The logarithmic factor in the update time comes from the fact, that one has to insert the items to the designated position since every list in the data structure is sorted. The number of items that we consider during an update step is equal to the number of items we would consider if the data structure was not lexicographically ordered. The logarithmic factor in the update step of Lemma 4.31 comes from the fact that we need logarithmic time to insert or delete an element to/from the corresponding sum tree. This has to be done for every item whose fit status changes during an update step. If the data structure is lexicographically ordered this has also to be done for the same number of items as in the non-lexicographically ordered structure. Thus, it follows, when considering a lexicographically ordered version of the data structure of Lemma 8.1, we obtain the following result.

**Corollary 8.4.** *There is a dynamic algorithm that receives a  $q$ -hierarchical  $k$ -ary CQ (without aggregates)  $Q(x_1, \dots, x_k, b_1, \dots, b_\ell)$  and a  $\sigma$ -db  $D_0$  of size  $\|D_0\|$ , and computes within preprocessing time  $\text{poly}(Q) \cdot O(\|D_0\| \log(\|D_0\|))$  a data structure that can be updated in time  $\text{poly}(Q) \cdot O(\log(\|D\|))$  and allows the following:*

- (a) *output the query result size  $|Q(D)|$  in time  $O(1)$ ,*
- (b) *enumerate the tuples in  $Q(D)$  in lexicographical order with delay  $\text{poly}(Q)$ ,*
- (c) *test for an input tuple  $\bar{a} \in \text{dom}^{k+\ell}$  if  $\bar{a} \in Q(D)$  within time  $t_t = \text{poly}(Q)$ , and*

8. Outputting the  $j$ th solution for  $q$ -hierarchical conjunctive queries under updates

(d) upon input of a tuple  $\bar{b} \in \mathbf{dom}^{k+\ell}$  output the tuple

$$\max \{ \bar{a} \in Q(D) : \bar{a} \leq \bar{b} \}$$

if it exists, or a *SmallerThanMinimum*-message otherwise in time  $t_i = \text{poly}(Q)$ ,

(e) upon input of an arbitrary number  $j$ , take time  $\text{poly}(Q) \cdot O(\log(\|D\|))$  to immediately output the  $j$ th tuple that the enumeration procedure of (b) would output.

(f) upon input of an arbitrary tuple  $\bar{a} \in Q(D)$ , take time  $\text{poly}(Q) \cdot O(\log(\|D\|))$  to immediately output the number  $j$  such that  $\bar{a}$  is the  $j$ th tuple the enumeration procedure of (b) would output,

where  $D$  is the current database.

Theorem 3.9 is the result of applying Lemma 6.18 on Corollary 8.4:

Suppose that we receive as input a  $q$ -hierarchical CQ with aggregates  $Q$ . To initialise the data structure for  $Q$  on database  $D$  the algorithm we obtain from Lemma 6.18(b) to compute a corresponding schema  $\tilde{\sigma}$  and a  $q$ -hierarchical query  $\tilde{Q}$  (without aggregates). Using Corollary 8.4, we obtain that there is a data structure for  $\tilde{Q}$  and  $\tilde{\sigma}$ -dbs that can be initialised for the empty database in time  $\tilde{t}_i = \text{poly}(Q)$  and can be updated in time  $\tilde{t}_u = \text{poly}(\tilde{Q}) \cdot O(\log(\|\tilde{D}\|)) \leq \text{poly}(Q) \cdot O(\log(\|D\|))$ . Using Lemma 6.18(c) we obtain an algorithm for  $Q$  with initialisation time  $t_i = \text{poly}(Q)$  and update time  $t_u = \text{poly}(Q) \cdot O(\log(\|D\|))t_a$ . Note that the data structure from Corollary 8.4 for  $\tilde{Q}$  on  $\tilde{D}$  is present.

The **count**-routine can be done as follows. Use the **count**-routine from Corollary 8.4(a) for  $\tilde{Q}$  and  $\tilde{D}$ . Since  $Q(D) = \tilde{Q}(\tilde{D})$  (Lemma 6.18(a)), it follows that we output the number  $\tilde{Q}(\tilde{D})$ .

The **enumerate**-routine can be done as follows. Use the **enumerate**-routine from Corollary 8.4(b) for  $\tilde{Q}$  and  $\tilde{D}$ . Since  $Q(D) = \tilde{Q}(\tilde{D})$  (Lemma 6.18(a)), it follows that we enumerate the tuples in  $Q(D)$ . Moreover, the tuples will be enumerated lexicographically.

The **test**-routine can be done as follows. Use the **test**-routine from Corollary 8.4(b) for  $\tilde{Q}$  and  $\tilde{D}$ . Since  $Q(D) = \tilde{Q}(\tilde{D})$  (Lemma 6.18(a)), it follows that we test for the tuples in  $Q(D)$  upon input of an tuple.

The routine in Theorem 3.9(d) for computing the next smaller tuple can be done as follows. Use the corresponding routine from Corollary 8.4(e). Since  $Q(D) = \tilde{Q}(\tilde{D})$  (Lemma 6.18(a)), it follows that we obtain the next smaller one from  $Q(D)$ .

The **jth**-routine can be done as follows. Use the **jth**-routine from Corollary 8.4(e). Since we use the **enumerate** routine from Corollary 8.4(b) it follows that we output the  $j$ th tuple the enumeration routine would output.

The **jth**-reverse-routine can be done as follows. Use the **jth**-reverse-routine from Corollary 8.4(f). Since we use the **enumerate** routine from Corollary 8.4(b) it follows that we output the  $j$ th tuple the enumeration routine would output.

## 8.5. Lower bounds

This section is devoted to the proof of Theorem 3.16.

For proving the theorem the following lemma will be convenient:

**Lemma 8.5.** *Let  $Q$  be a CQ and let  $D$  be a database and  $\leq$  be a linear order on the result set  $Q(D)$ . Suppose that there is an Algorithm  $A_{\leq}$  that outputs in time  $t_j$  on input of a number  $j \in \mathbb{N}$  the tuple  $\bar{a} \in Q(D)$  such that  $\bar{a}$  is the  $j$ th tuple with respect to  $\leq$  if  $j \leq |Q(D)|$  or outputs the **OutOfRange**-Message otherwise. Then, there is an algorithm  $B$  that enumerates  $Q(D)$  with delay  $t_d = t_j$ .*

*Proof.* See Algorithm 14 for Algorithm B.

---

**Algorithm 14** Algorithm B.

---

```

1:  $j = 1$ 
2: Let  $\bar{a}$  be the output of  $A_{\leq}$  on input  $j$ .
3: while  $\bar{a}$  is not the OutOfRange-Message do
4:   Output  $\bar{a}$ .
5:    $j \leftarrow j + 1$ 
6:   Let  $\bar{a}$  be the output of  $A_{\leq}$  on input  $j$ .
7: Output end-of-enumeration message EOE.

```

---

Clearly, B enumerates  $Q(D)$  with delay  $t_j$ . □

*Proof of Theorem 3.16.* Suppose for a contradiction that there is an algorithm that receives a self join free conjunctive query  $Q$  that is not q-hierarchical and a database  $D$ , and computes with arbitrary preprocessing time a data structure that can be updated in time  $O(n^{1-\varepsilon})$  and outputs in time  $O(n^{1-\varepsilon})$  on input of a  $j \in \mathbb{N}$  the tuple  $\bar{a} \in Q(D)$  such that  $\bar{a}$  is the  $j$ th tuple with respect to a linear order  $\leq$  on  $Q(D)$  if  $j \leq |Q(D)|$  or outputs the **OutOfRange**-Message otherwise for any linear order of the result set  $Q(D)$ . Applying Lemma 8.5, we can enumerate the tuples in  $Q(D)$  with delay  $O(n^{1-\varepsilon})$ . This violates to Theorem 3.13. □

This concludes the proof of Theorem 3.16.



## 9. Answering unions of conjunctive queries under updates

In this chapter we consider dynamic query evaluation for UCQs. The syntax and semantics of UCQs are given in page 20. To transfer our notions of *hierarchical* queries from CQs to UCQs, we say that a UCQ  $Q(\bar{u})$  of the form  $Q_1(\bar{u}_1) \cup \dots \cup Q_d(\bar{u}_d)$  is q-hierarchical (t-hierarchical) if every CQ  $Q_i(\bar{u}_i)$  in the union is q-hierarchical (t-hierarchical). Note that for *Boolean* queries (CQs as well as UCQs) the notions of being q-hierarchical and being t-hierarchical coincide, and for a  $k$ -ary UCQ  $Q$  it can be checked in time  $\text{poly}(Q)$  if  $Q$  is q-hierarchical or t-hierarchical.

In the next section we consider the task of testing and enumerating UCQs. In Section 9.2 we identify the UCQs for which the task of reporting the  $j$ th query result can be done efficiently under updates.

### 9.1. Enumerating and testing UCQs

**Testing.** This paragraph is based on [21]. Theorem 3.10 generalises the statement of Theorem 3.6 from CQs to UCQs. Its proof follows easily from the Theorems 3.2 and 3.6. Theorem 3.10 is restated in the following.

**Theorem 3.10** ([21]). *There is a dynamic algorithm that receives a  $t$ -hierarchical  $k$ -ary UCQ  $Q$  and a  $\sigma$ -db  $D_0$ , and computes within  $t_p = \text{poly}(Q) \cdot O(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)$  and allows to test for an input tuple  $\bar{a} \in \text{dom}^k k$  if  $\bar{a} \in Q(D)$  within time  $t_t = \text{poly}(Q)$ . Furthermore, the algorithm allows to answer a  $t$ -hierarchical Boolean UCQ within time  $t_{ans} = O(1)$ .*

*Proof of Theorem 3.10.* It follows immediately from Theorem 3.6 (and Theorem 3.2 for the statement on *Boolean* UCQs), as we can maintain all CQs in the union in parallel and then decide whether at least one of them is satisfied by the current database and the given tuple.  $\square$

**Enumerating.** It turns out that q-hierarchical UCQs, like q-hierarchical CQs, allow for efficient enumeration under updates. This is stated in Theorem 3.11 restated in the following.

**Theorem 3.11** ([21]). *There is a dynamic algorithm that receives a  $q$ -hierarchical  $k$ -ary UCQ  $Q$  and a  $\sigma$ -db  $D_0$ , and computes within  $t_p = \text{poly}(Q) \cdot O(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u = \text{poly}(Q)$  and allows to enumerate  $Q(D)$  with delay  $t_d = \text{poly}(Q)$ .*

## 9. Answering unions of conjunctive queries under updates

In contrast to Theorem 3.10, the result does not follow immediately from the tractability of the enumeration problem for q-hierarchical CQs, because one has to ensure that tuples from result sets of two different CQs are not reported twice while enumerating their union.

In the following an alternative proof of Theorem 3.11 as in [21] is given.

To prove Theorem 3.11, we first show a general method for enumerating the union of sets which is stated in Lemma 9.1. This idea relies on the algorithm for enumerating the union of two sets in [39].

**Lemma 9.1.** *Let  $S_1, \dots, S_m$  be  $m$  sets. Suppose that there is for every  $i \in [m]$  an algorithm  $A_{\text{enum}}^i$  that enumerates the elements in  $S_i$  with delay  $t_d$ , and an algorithm  $A_{\text{test}}^i$  that can test in time  $t_t$  for a given element  $a$  whether  $a \in S_i$ . Then, there is an algorithm that enumerates the elements in  $\bigcup_{i \in [m]} S_i$  without repetition with delay  $O(n(t_t + t_d))$ .*

To illustrate the idea of Lemma 9.1 we consider the following example.

**Example 9.2.** *Let us consider the three sets  $S_1 := \{1, 4, 5\}$ ,  $S_2 := \{2, 4\}$ ,  $S_3 := \{2, 3, 4, 5\}$  and we have for each  $i \in [3]$  an algorithm  $A_{\text{enum}}^i$  that enumerates the set  $S_i$  with delay  $t_d$  and an algorithm  $A_{\text{test}}^i$  which upon input of a natural number  $n$  tests if  $n \in S_i$ . The aim is to enumerate  $S_1 \cup S_2 \cup S_3$  without repetition. The algorithm for the general case is given in Algorithm 15. First of all, we set for all  $i \in [3]$  the element  $s_i$  to the first element of  $S_i$ , i.e.  $s_1 = 1, s_2 = 2, s_3 = 2$ . Then, we start with the set  $S_1$  (with  $\ell = 1$  in the for-loop in line 2) and start with OUTPUT(1). In the procedure we check in the for-loop if there is a  $K \in \{2, 3\}$  with  $1 \in S_K$ , i.e., if  $1 \in S_2 \cup S_3$ . Since this is not the case, we print 1 and set  $s_1$  to 4 and since  $s_1 \neq \text{EOE}$  we again execute OUTPUT(1). Since  $4 \in S_2$  the if-condition in line 9 holds for  $K = 2$ . Thus, we continue with setting  $s_1$  to 5 and start with OUTPUT(2). Since  $s_2 = 2 \in S_3$ , the if condition in line 9 holds for  $K = 3$  and set  $s_2$  to 4 and continue with OUTPUT(3). This procedure prints  $s_3 = 2$  and sets  $s_3$  to 3. Then, we continue with the while-loop since  $s_1 \neq \text{EOE}$  and since  $5 \in S_3$  we set  $s_1$  to EOE and execute OUTPUT(3) which prints 3 and sets  $s_3$  to 4. Since  $s_1 = \text{EOE}$ , the for-loop in line 2 sets  $\ell$  to 2 and thus, we consider the set  $S_2$ . Since  $4 \in S_3$ , we print  $s_3 = 4$  and set  $s_3$  to the next element in  $S_3$ , which is 5 and  $s_2$  gets the value EOE and afterwards we enumerate 5 from  $S_3$  and then the algorithm terminates since all elements in  $S_3$  were printed and  $s_3 = \text{EOE}$ .*

*The elements that were printed in OUTPUT(1) belong to  $S_1 \setminus (S_2 \cup S_3)$ , the elements that were printed in OUTPUT(2) belong to  $S_2 \setminus S_3$  and the elements in OUTPUT(3), that were printed, belong to  $S_3$ . Altogether, we enumerated the set  $S_1 \cup S_2 \cup S_3$ .*

We consider now the general case.

*Proof of Lemma 9.1.* Algorithm 15 shows how to enumerate the set  $S_1 \cup S_2 \cup \dots \cup S_m$ .

**Correctness** Every time we call OUTPUT( $\ell$ ), we print the element  $s_\ell$  if and only if there is no  $K \in \{\ell + 1, \dots, m\}$  such that  $s_\ell \in S_K$ , i.e.,  $s_\ell \notin \bigcup_{K=\ell+1}^m S_K$ . The fact

---

**Algorithm 15** Algorithm for enumerating  $S_1 \cup S_2 \cup \dots \cup S_m$ .

---

```

1: For all  $i \in [m]$  let  $s_i$  be the first element of  $S_i$ .
2: for  $\ell = 1$  to  $m$  do
3:   while  $s_\ell \neq \text{EOE}$  do
4:     OUTPUT( $\ell$ )
5: print EOE
6:
7: procedure OUTPUT( $\ell$ )
8:   for  $K = \ell + 1$  to  $m$  do
9:     if  $s_\ell \in S_K$  then
10:      Set  $s_\ell$  to the successor element of  $s_\ell$  in  $S_\ell$ 
11:      OUTPUT( $K$ )
12:   return
13: print  $s_\ell$ 
14: Set  $s_\ell$  to the successor element in  $S_\ell$ 

```

---

that we initialise  $s_\ell$  to the first element and then we set  $s_\ell$  to its successor element in every call of OUTPUT( $\ell$ ) and that the for and while loop at the beginning ensures that we repeat the call OUTPUT( $\ell$ ) until  $s_\ell = \text{EOE}$  for all  $\ell \in [m]$ , guarantees that we consider every element in  $S_\ell$  exactly once. All in all, the elements that are printed in the whole algorithm are

$$S_1 \setminus \left( \bigcup_{K=2}^m S_K \right) \cup \dots \cup S_m = \bigcup_{\ell=1}^m \left[ S_\ell \setminus \left( \bigcup_{K=\ell+1}^m S_K \right) \right] = \bigcup_{\ell=1}^m S_\ell$$

To show that we enumerate the elements without repetition, let us assume for a contradiction that there are  $\ell_1, \ell_2 \in [m]$  with  $\ell_1 < \ell_2$  and an element  $s \in S_{\ell_1} \cap S_{\ell_2}$  such that  $s$  will be printed during a call of OUTPUT( $\ell_1$ ) and a call of OUTPUT( $\ell_2$ ). Since  $s \in S_{\ell_2}$  and  $\ell_2 > \ell_1$  this will not be print by a call of OUTPUT( $\ell_1$ ). Furthermore, since we enumerate the sets themselves without repetition, we enumerate the disjunction without repetition.

**Running time** First of all, we show that every time we call OUTPUT( $\ell$ ) it holds that  $s_\ell \neq \text{EOE}$ . If OUTPUT( $\ell$ ) will be called in line 4 this is guaranteed by the condition of the while loop, and if it is called in line 11 there is a  $k < \ell$  and an element  $s \in S_k$  such that  $m \in S_\ell$  and  $s \notin S_{k+1} \cup \dots \cup S_{\ell-1}$  (since otherwise the algorithm calls OUTPUT( $p$ ) [instead of OUTPUT( $\ell$ )] for the smallest  $p \in \{k+1, \dots, \ell-1\}$  with  $m \in S_p$ ). The number of times that the procedure OUTPUT( $\ell$ ) will be called in line 11 is at most  $|(S_1 \cup \dots \cup S_{\ell-1}) \cap S_\ell| \leq |S_\ell|$  (since there must be an element  $s_\ell \in S_\ell$  and a  $k \in [\ell-1]$  such that  $s_\ell \in S_k \setminus (S_{k+1} \cup \dots \cup S_{\ell-1})$ ). In particular, every element in  $S_\ell$  will be mentioned in  $S_1 \cup \dots \cup S_{\ell-1}$  once). Therefore  $s_\ell \neq \text{EOE}$  if it is called in line 11. When we call OUTPUT( $\ell$ ), we either output an element or there is a  $K > \ell$  such that OUTPUT( $K$ ) will be called. Since  $K \leq m$ , we need at most  $m$  steps until we output

## 9. Answering unions of conjunctive queries under updates

an element, after a call of `OUTPUT`( $\ell$ ). During the call, we have at most  $K - \ell$  tests in time  $O(t_t)$  and we jump to the next element in a set in time  $O(t_d)$ . Therefore, it takes time  $O(m(t_d + t_t))$  after a call of the `OUTPUT` procedure until we output an element. By construction of the algorithm, we have  $O(m(t_d + t_t))$  delay.  $\square$

Using Lemma 9.1 we obtain a proof for Theorem 3.11:

*Proof of Theorem 3.11.* Let  $Q = Q_1 \cup \dots \cup Q_m$  be the input query. Note that  $\|Q_i\| \leq \|Q\|$ . We store in parallel for every  $i \in [m]$  the data structure from Theorem 3.3 for the q-hierarchical query  $Q_i$  on  $D_0$ . Every time the database receives an update we update the data structure for every  $i \in [m]$  using the update procedure from Theorem 3.3. This takes time  $O(m \cdot \text{poly}(Q)) = O(\text{poly}(Q))$ . To enumerate the result sets  $Q_i(D)$  use the enumeration algorithm from Theorem 3.3(b) to enumerate the tuples with delay  $t_d = O(\text{poly}(Q))$ . Furthermore, with Theorem 3.3(a) we can test time  $t_t = O(\text{poly}(Q))$ , whether a tuple belongs to the result set  $Q_i(D)$ . Applying the enumeration and testing routines to Lemma 9.1 we obtain an enumeration algorithm for  $Q(D)$  that enumerates the tuples without repetition and with delay  $O(m(\text{poly}(Q) + \text{poly}(Q))) = O(\text{poly}(Q))$ .

This concludes the proof of Theorem 3.11.  $\square$

## 9.2. Reporting the $j$ th solution

In this section we identify a class of UCQs for which we can maintain the  $j$ th solution problem under updates. Here we consider strongly exhaustively q-hierarchical queries.

Recall that a union of conjunctive queries  $\bigcup_{i \in [m]} Q_i$  is strongly exhaustively q-hierarchical if there is a q-hierarchical query  $Q'$  with  $Q'(D) = \bigcap_{i \in [m]} Q_i(D)$  for all  $\sigma$ -dbs  $D$ .

Let  $Q = Q_1 \cup \dots \cup Q_m$  be an arbitrary strongly exhaustively q-hierarchical union of conjunctive queries.

In the following Theorem 3.12 is restated.

**Theorem 3.12.** *There is a dynamic algorithm that receives a strongly exhaustively q-hierarchical  $k$ -ary UCQ  $Q$  and a  $\sigma$ -db  $D_0$  of size  $\|D_0\|$ , and computes within preprocessing time  $\exp(Q) \cdot O(\|D_0\| \log(\|D_0\|))$  a data structure that can be updated in time  $\exp(Q) \cdot O(\log(\|D\|))$  and allows the following:*

- (a) *enumerate the tuples in  $Q(D)$  with delay  $\text{poly}(Q)$  and*
- (b) *upon input of an arbitrary number  $j$ , take time  $\exp(Q) \cdot O(\log(\|D\|))$  to immediately output the  $j$ th tuple that the enumeration procedure of (a) would output,*

*where  $D$  is the current database.*

To prove Theorem 3.12 we use the following lemma that enables to output the  $j$ th solution when using the algorithm from Lemma 9.1.

**Lemma 9.3.** *Let  $X$  be a set and let  $S_1, \dots, S_m$  be  $m$  subsets of  $X$  and for all  $\emptyset \neq \mathcal{I} \subseteq [m]$  let  $S_{\mathcal{I}} := \bigcap_{i \in \mathcal{I}} S_i$ . Assume that for all  $\emptyset \neq \mathcal{I} \subseteq [m]$  we have the following algorithms available.*



## 9.2. Reporting the $j$ th solution

- (a) An algorithm  $A_{\text{enum}}^{\mathcal{I}}$  that enumerates elements in  $S_{\mathcal{I}}$  with delay  $t_d$  in lexicographical order and
- (b) an algorithm  $A_{\text{count}}^{\mathcal{I}}$  that outputs the number  $|S_{\mathcal{I}}|$  in time  $t_c$  and
- (c) an algorithm  $A_{\text{test}}^{\mathcal{I}}$  that upon input of an element  $x \in X$  tests in time  $t_t$  if  $x \in S_{\mathcal{I}}$  and
- (d) an algorithm  $A_{\text{jth}}^{\mathcal{I}}$  that upon input of a number  $j \in \mathbb{N}$  returns the  $j$ th element in  $S_{\mathcal{I}}$  that the algorithm  $A_{\text{enum}}^{\mathcal{I}}$  will output in time  $t_j$  if it exists or a **OutOfRange**-message otherwise and
- (e) an algorithm  $A_{\text{jth-rev}}^{\mathcal{I}}$  that upon input of an element  $a \in S_{\mathcal{I}}$  outputs the number  $j$  in time  $t_{jr}$  such that  $a$  is the  $j$ th element algorithm  $A_{\text{enum}}^{\mathcal{I}}$  will output and
- (f) an algorithm  $A_{\text{ex}}^{\mathcal{I}}$  that upon input of an element  $x \in X$  outputs in time  $t_l$  the element  $\max\{a : a \in S_{\mathcal{I}}, a \leq x\}$  if it exists or a **SmallerThanMinimum** message otherwise.

Then, there is an algorithm that upon input of a number  $j \in \mathbb{N}$  outputs the  $j$ th element that the algorithm in Lemma 9.1 for enumerating  $S_1 \cup \dots \cup S_m$ , by applying the enumeration algorithms  $A_{\text{enum}}^{\{i\}}$  and testing algorithms  $A_{\text{test}}^{\{i\}}$  for all  $i \in [m]$ , will output if  $j \leq |S_1 \cup \dots \cup S_m|$  or output an **OutOfRange**-message otherwise. The algorithm takes time  $O(m^2 2^m (t_j + t_l + t_{jr} + t_c) + n(t_t + t_j))$ .

The  $j$ th step is the fragment of Algorithm 15 that enumerates the  $j$ th element. To be precise, the  $j$ th step starts, when we start the **OUTPUT**( $\ell$ ) procedure in line 4 in Algorithm 15 that ends by printing the  $j$ th element.

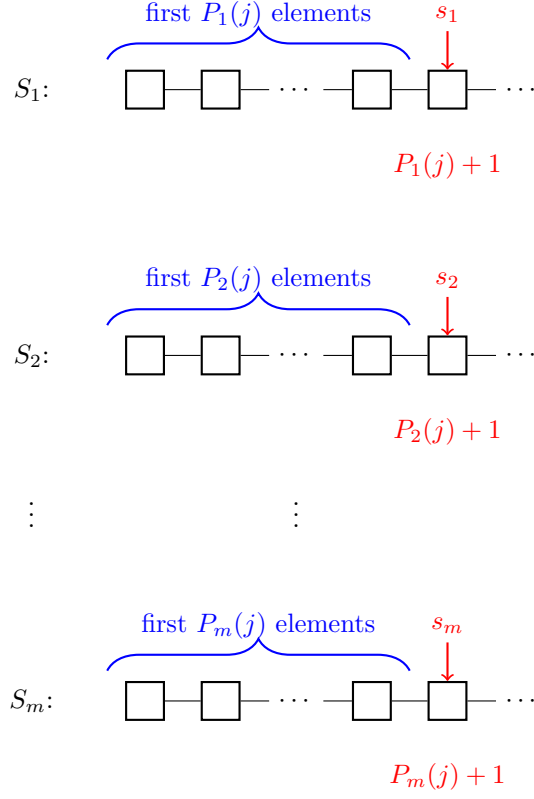
To prove Lemma 9.3, we try to simulate the  $j$ th step of Algorithm 15, i.e., the fragment of the algorithm that enumerates the  $j$ th element. To do this, we need to compute the current elements  $s_i$  for all  $i \in [m]$  and the current number  $\ell \in [m]$  in Algorithm 15 before the  $j$ th step starts.

To compute these elements we define for all  $k \in [m]$  and  $j \in \mathbb{N}$  numbers  $P_k(j)$  such that the following holds. The elements  $s_k$  (this is the current element of the set  $S_k$  in Algorithm 15) is set to the  $(P_k(j)+1)$ th element in  $S_k$  for all  $k \in [m]$  when the  $j$ th step starts. In particular, the number of elements in  $S_k$  we have already considered, i.e., printed or skipped is  $P_k(j)$ . Note that the element  $s_k$ , i.e., the  $(P_k(j)+1)$ th element of  $S_k$ , was neither printed nor skipped. See Figure 9.1 for an illustration of the numbers  $P_k(j)$  when the  $j$ th step starts. We will later describe how to compute the number  $P_k(j)$ . To determine the number  $\ell$  in the for-loop in line 2 in Algorithm 15 we use the number  $j$  in the following way. For all  $i \in [j]$  one has to check if  $|S_1 \cup \dots \cup S_{i-1}| < j \leq |S_1 \cup \dots \cup S_i|$ . The number  $i$  for which this condition holds true is the sought number for  $\ell$ . The set  $S_\ell$  is the current set which is considered since already  $|S_1 \cup \dots \cup S_{\ell-1}|$  elements were printed by the algorithm and, in particular, all elements in  $S_1 \cup \dots \cup S_{\ell-1}$  were either printed or skipped by the algorithm, but not all elements in  $S_\ell$ .

The algorithm that simulates the  $j$ th step of Algorithm 15 is given in Algorithm 16. We will later describe how to implement the lines of Algorithm 15.

## 9. Answering unions of conjunctive queries under updates

Figure 9.1.: Illustration of how the numbers  $P_k(j)$  are set at the beginning of the  $j$ th step (for the special case where  $\ell = 1$ ).




---

**Algorithm 16** Algorithm to output the  $j$ th element of  $S_1 \cup \dots \cup S_m$ .

---

- 1: **Input:**  $j \in \mathbb{N}$ .
  - 2: **if**  $j > |S_1 \cup \dots \cup S_m|$  **then**
  - 3:     Output **OutOfRange**-message
  - 4:   Let  $\ell$  be the number such that  $|S_1 \cup \dots \cup S_{\ell-1}| < j \leq |S_1 \cup \dots \cup S_\ell|$ .
  - 5:   Let  $a$  be the  $(P_\ell(j)+1)$ th element of  $S_\ell$ .
  - 6:   Let  $k := \ell$ .
  - 7:   **while** there is a  $K > k$  such that  $a \in S_K$  **do**
  - 8:      $k := \min \{K : K > k \text{ and } a \in S_K\}$
  - 9:     Let  $a$  be the  $(P_k(j)+1)$ th element of  $S_k$ .
  - 10: Output  $a$
-

## 9.2. Reporting the $j$ th solution

In the  $j$ th step, we iterate over the elements in  $S_\ell$  and therefore we consider the  $(P_\ell(j)+1)$ th element of  $S_\ell$ . (Note that by definition of  $\ell$  we have  $P_\ell(j) + 1 \leq |S_1 \cup \dots \cup S_\ell|$ .) This is the  $\ell$  of the first for-loop in Algorithm 15 when enumerating the  $j$ th element. Then, the enumeration algorithm checks (in the OUTPUT procedure in Algorithm 15) for  $k \in \{\ell + 1, \dots, m\}$  in ascending order if  $a \in S_k$  holds. If we find such a  $k$ , we jump to the next element in  $S_k$  that was not output or skipped by the enumeration algorithm. This is exactly the  $(P_k(j)+1)$ th element of  $S_k$ . The enumeration algorithm repeats this until we cannot find an element that appears in a set with a higher index. This is exactly the element that is output by Algorithm 16 upon input of the number  $j$ .

The hard part is to compute the numbers  $P_k(j)$ . We will use the following notation. For all  $k \in [m]$  let  $a_{P_k(j)}$  be the  $P_k(j)$ th element in the enumeration of  $S_k$ . Furthermore, let

$$\tilde{S}_k := \begin{cases} \{s \in S_k : s \leq a_{P_k(j)}\} & \text{if } k \geq \ell \\ S_k & \text{otherwise} \end{cases}.$$

In other words:  $\tilde{S}_k$  is the set of the elements in  $S_k$  that Algorithm 15 already printed or skipped until it started the  $j$ th step. Before we go into detail for the general case let us consider the following example.

**Example 9.4.** Recall Example 9.2. Our aim is now to output the 3rd element. Algorithm 16 sets  $\ell = 1$  and since we iterate over the elements in  $S_1$ , we have to consider the 3rd element in  $S_1$ , i.e.,  $P_1(j) = j = 3$ . We consider 5 which is the 3rd element in  $S_1$  and since  $5 \in S_3$  we jump to the set  $S_3$ . Now we count how often we jumped to Output(3) during the enumeration of the first two elements. There are two cases why this could happen. The first case is that there is an element  $a$  in  $\tilde{S}_1$  (one of the elements in  $S_1$  we already considered) which appears in  $S_3$  but not in  $S_2$ . Note that if  $a$  appeared in  $S_2$  we would first jump to OUTPUT(2) and not to  $S_3$ . This is exactly the case when considering  $5 \in S_1$  in Example 9.2. The case can be written as follows. There is an element  $a \in (\tilde{S}_1 \setminus S_2) \cap S_3$ . The second case is that there is an element in  $S_2$ , which we already considered, that also appears in  $S_3$ , i.e., there is an element in the set  $\tilde{S}_2 \cap S_3$ . All in all, we can express the number  $P_3(j)$  as

$$\begin{aligned} P_3(j) &= \left| ((\tilde{S}_1 \setminus S_2) \cap S_3) \cup (\tilde{S}_2 \cap S_3) \right| \\ &\stackrel{(\star)}{=} \left| (\tilde{S}_1 \setminus S_2) \cap S_3 \right| + \left| \tilde{S}_2 \cap S_3 \right| \\ &\stackrel{(\star\star)}{=} \left| (\tilde{S}_1 \cap S_3) \setminus (\tilde{S}_1 \cap S_2 \cap S_3) \right| + \left| \tilde{S}_2 \cap S_3 \right| \\ &\stackrel{(\star\star\star)}{=} \left| \tilde{S}_1 \cap S_3 \right| - \left| \tilde{S}_1 \cap S_2 \cap S_3 \right| + \left| \tilde{S}_2 \cap S_3 \right| \end{aligned}$$

Equation  $(\star)$  follows from the fact that the unions are disjoint, equation  $(\star\star)$  by the de Morgan's law and  $(\star\star\star)$  from the fact that  $(\tilde{S}_1 \cap S_2 \cap S_3)$  is a subset of  $(\tilde{S}_1 \cap S_3)$ .

## 9. Answering unions of conjunctive queries under updates

It remains to show how to compute these cardinalities. In the following we show how to compute the number  $\left|(\tilde{S}_1 \cap S_2 \cap S_3)\right|$ . Suppose that we know the number  $P_1(j)$ . Then, we use algorithm  $A_{jth}^{\{1\}}$  to compute the  $P_1(j)$ th element  $a$  of  $S_1$ . Then, using algorithm  $A_{lex}^{\{1,2,3\}}$  we compute the maximum element  $b \in (\tilde{S}_1 \cap S_2 \cap S_3)$  which is lexicographically smaller than or equal to  $a$ . If such an element does not exist,  $\left|(\tilde{S}_1 \cap S_2 \cap S_3)\right| = 0$ . Otherwise, we obtain the number  $\left|(\tilde{S}_1 \cap S_2 \cap S_3)\right|$  by applying  $A_{jth-rev}^{\{1,2,3\}}$  upon input of  $b$ . This is the number of elements in  $S_1 \cap S_2 \cap S_3$  that are smaller than the current  $s_1$  value (the  $P_1(j)$ th element in  $S_1$ ). The number  $\left|\tilde{S}_2 \cap S_3\right|$  can be computed by a similar way and the cardinality of  $S_2 \cap S_3$  can simply be determined by using the  $A_{count}^{\{2,3\}}$  algorithm.

We consider now the general case. To compute the number  $P_k(j)$  it is necessary to have the  $P_i(j)$  value available for all  $i < k$ . Hence, we compute the number  $P_i(j)$  in ascending order, i.e., for  $i = 1, 2, \dots, k$ . In the following it is described how to obtain a formula to compute the number  $P_j(k)$  for  $k > \ell$ . We have to count the number of elements in  $S_k$  we have jumped to, when considering the previous elements. An element in  $S_k$  will only be considered if there is an  $u \leq k$  such that there is an element  $a \in S_u$  which also appears in  $S_k$  but not in  $S_{u+1} \cup \dots \cup S_{k-1}$  (since otherwise the OUTPUT operation would jump to the corresponding set  $S_q$  with  $q \in \{u+1, \dots, k-1\}$ ). In other words, there is an  $i < k$  such that there is an element  $a \in (S_u \setminus \bigcup_{j=u+1}^{k-1} S_j) \cap S_k$ . To describe the number  $P_k(j)$  one has to count the number of element for which the conditions above holds. Therefore, we obtain for  $k > \ell$  the following

$$P_k(j) = \left| \bigcup_{u=1}^{k-1} \left( \tilde{S}_u \setminus \bigcup_{i=u+1}^{k-1} S_i \right) \cap S_k \right|. \quad (9.1)$$

By using some transformations of the formula above, we obtain the following for-

mula.

$$P_k(j) = \left| \bigcup_{u=1}^{k-1} \left( \tilde{S}_u \setminus \bigcup_{i=u+1}^{k-1} S_i \right) \cap S_k \right| \quad (9.2)$$

$$= \left| \bigcup_{u=1}^{k-1} \left( (\tilde{S}_u \cap S_k) \setminus \left( \bigcup_{i=u+1}^{k-1} S_i \cap \tilde{S}_u \cap S_k \right) \right) \right| \quad (9.3)$$

$$\stackrel{(\star)}{=} \sum_{u=1}^{k-1} \left| (\tilde{S}_u \cap S_k) \setminus \left( \bigcup_{i=u+1}^{k-1} S_i \cap \tilde{S}_u \cap S_k \right) \right| \quad (9.4)$$

$$\stackrel{(\star\star)}{=} \sum_{u=1}^{k-1} \left( |\tilde{S}_u \cap S_k| - \left| \bigcup_{i=u+1}^{k-1} S_i \cap \tilde{S}_u \cap S_k \right| \right) \quad (9.5)$$

$$\stackrel{(\star\star\star)}{=} \sum_{u=1}^{k-1} \left( |\tilde{S}_u \cap S_k| - \sum_{\emptyset \neq \mathcal{I} \subseteq \{u+1, \dots, k-1\}} (-1)^{|\mathcal{I}|-1} \left| \bigcap_{i \in \mathcal{I}} S_i \cap \tilde{S}_u \cap S_k \right| \right) \quad (9.6)$$

$$= \sum_{u=1}^{k-1} \left( |\tilde{S}_u \cap S_k| + \sum_{\emptyset \neq \mathcal{I} \subseteq [k-1] \setminus [u]} (-1)^{|\mathcal{I}|} \left| \bigcap_{i \in \mathcal{I}} S_i \cap \tilde{S}_u \cap S_k \right| \right) \quad (9.7)$$

The equation  $(\star)$  follows from the fact that the unions are disjoint and  $(\star\star)$  from the fact that set  $\bigcup_{i=u+1}^{k-1} S_i \cap \tilde{S}_u \cap S_k$  is a subset of  $(\tilde{S}_u \cap S_k)$ . In  $(\star\star\star)$  we used the inclusion-exclusion principle (see Equation 2.2).

**Claim 9.5.** *Let  $\ell$  be the number such that  $|S_1 \cup \dots \cup S_{\ell-1}| < j \leq |S_1 \cup \dots \cup S_\ell|$ . Then, the following holds.*

•

$$P_\ell(j) = j - \sum_{\emptyset \neq \mathcal{I} \subseteq [\ell]} (-1)^{|\mathcal{I}|+1} \left| \bigcap_{i \in \mathcal{I}} S_i \right| \quad (9.8)$$

and

• for all  $k > \ell$  we have

$$P_k(j) = \sum_{u=1}^{k-1} \left( |\tilde{S}_u \cap S_k| + \sum_{\emptyset \neq \mathcal{I} \subseteq [k-1] \setminus [u]} (-1)^{|\mathcal{I}|} \left| \bigcap_{i \in \mathcal{I}} S_i \cap \tilde{S}_u \cap S_k \right| \right). \quad (9.9)$$

*Proof.* The case for  $P_k(j)$  with  $k \in [m] \setminus [\ell]$  follows immediately from Equation (9.7).

$S_\ell$  is the set which is iterated during the  $j$ th step. The number  $j - |S_1 \cup \dots \cup S_{\ell-1}|$  is the number of elements we have to consider in  $S_\ell$  until we reach the  $j$ th element in  $S_1 \cup \dots \cup S_\ell$ , after the elements in  $|S_1 \cup \dots \cup S_{\ell-1}|$  were enumerated. The number of elements in  $S_\ell$  we considered until we enumerated the elements in  $S_1 \cup \dots \cup S_{\ell-1}$  is  $\left| \left( \bigcup_{u=1}^{\ell-1} S_u \setminus \bigcup_{i=u+1}^{\ell-1} S_i \right) \cap S_\ell \right|$ .

### 9. Answering unions of conjunctive queries under updates

In summary, the number  $P_\ell(j)$  can be determined as follows.

$$\begin{aligned}
P_\ell(j) &= j - |S_1 \cup \dots \cup S_{\ell-1}| + \left| \left( \bigcup_{u=1}^{\ell-1} S_u \setminus \bigcup_{i=u+1}^{\ell-1} S_i \right) \cap S_\ell \right| \\
&= j - \sum_{\emptyset \neq \mathcal{I} \subseteq [\ell-1]} (-1)^{|\mathcal{I}|+1} \left| \bigcap_{i \in \mathcal{I}} S_i \right| + \left| \left( \bigcup_{u=1}^{\ell-1} [S_u \cap S_\ell] \setminus \bigcup_{i=u+1}^{\ell-1} [S_i \cap S_\ell \cap S_u] \right) \right| \\
&\stackrel{(\star)}{=} j - \sum_{\emptyset \neq \mathcal{I} \subseteq [\ell-1]} (-1)^{|\mathcal{I}|+1} \left| \bigcap_{i \in \mathcal{I}} S_i \right| + \sum_{u=1}^{\ell-1} \left| \left( [S_u \cap S_\ell] \setminus \bigcup_{i=u+1}^{\ell-1} [S_i \cap S_\ell \cap S_u] \right) \right| \\
&\stackrel{(\star\star)}{=} j - \sum_{\emptyset \neq \mathcal{I} \subseteq [\ell-1]} (-1)^{|\mathcal{I}|+1} \left| \bigcap_{i \in \mathcal{I}} S_i \right| + \sum_{u=1}^{\ell-1} |S_u \cap S_\ell| - \left| \bigcup_{i=u+1}^{\ell-1} [S_i \cap S_\ell \cap S_u] \right| \\
&\stackrel{(\star\star\star)}{=} j - \sum_{\emptyset \neq \mathcal{I} \subseteq [\ell-1]} (-1)^{|\mathcal{I}|+1} \left| \bigcap_{i \in \mathcal{I}} S_i \right| + \sum_{u=1}^{\ell-1} |S_u \cap S_\ell| \\
&\quad - \sum_{\emptyset \neq \mathcal{I} \subseteq [\ell-1]} (-1)^{|\mathcal{I}|+1} \left| \bigcap_{i \in \mathcal{I}} S_i \cap S_\ell \cap S_u \right|
\end{aligned}$$

The equation  $(\star)$  follows from the fact that the unions are disjoint and  $(\star\star)$  from the fact that set  $\bigcup_{i=u+1}^{\ell-1} S_i \cap S_\ell \cap S_u$  is a subset of  $(S_\ell \cap S_u)$ . In  $(\star\star\star)$  we used the inclusion-exclusion principle (see Equation 2.2).  $\square$

To avoid notational clutter, we define for every non-empty set  $\mathcal{I} \subseteq [m]$  the following set.

$$S_{\mathcal{I}} := \bigcap_{i \in \mathcal{I}} S_i$$

We obtain the numbers  $|S_{\mathcal{I}}|$  from  $\mathbf{A}_{\text{count}}^{\mathcal{I}}$  in time  $t_c$ . Since we compute the numbers  $|S_{\mathcal{I}}|$  for all  $\emptyset \neq \mathcal{I} \subseteq [\ell]$  it takes time  $2^\ell t_c$  to compute  $P_\ell(j)$ .

To compute  $P_k(j)$  for all  $k > \ell$  we use the formula from Claim 9.5. To compute the cardinality of sets of the form  $\left| \bigcap_{i \in \mathcal{I}} S_i \cap \tilde{S}_u \right|$  for  $u \in [m]$  and  $\mathcal{I} \subseteq [m]$  we do the following. If  $S_i$  is of the form  $\tilde{S}_i = S_i$ , we simply use algorithm  $\mathbf{A}_{\text{count}}^{\mathcal{I} \cup \{u\}}$ . In the other case, we use Algorithm 17.

Let  $a_1, \dots, a_{P_u(j)}$  be the first  $P_u(j)$  elements of  $S_u$  enumerated by  $\mathbf{A}_{\text{enum}}^{\{u\}}$ . In particular, we have  $a_1 < \dots < a_{P_u(j)}$ . Note that

$$\bigcap_{i \in \mathcal{I}} S_i \cap \tilde{S}_u = \{c \in \{a_1, \dots, a_{P_u(j)}\} : c \in S_{\mathcal{I}}\}$$

Clearly, by definition of  $S_{\mathcal{I}} \cap \tilde{S}_u$  all elements in the set are smaller than or equal to  $a_{P_u(j)}$ . Then, for  $b$  in line 2 of Algorithm 17 we have

$$b = \max \{c \in \{a_1, \dots, a_{P_u(j)}\} : c \in S_{\mathcal{I}}\}$$

---

**Algorithm 17** Algorithm for computing  $\left| \bigcap_{i \in \mathcal{I}} S_i \cap \tilde{S}_u \right|$ .

---

- 1: Use  $A_{jth}^{\{u\}}$  to get the  $P_u(j)$ th element of  $S_u$ . Let  $a_{P_u(j)}$  be the element.
  - 2: Use  $A_{lex}^{\mathcal{I} \cup \{u\}}$  to get  $b$ , the maximum element in  $S_{\mathcal{I} \cup \{u\}}$  which is smaller than (or equal to)  $a_{P_u(j)}$ .
  - 3: **if**  $b == \text{SmallerThanMinimum}$  **then**
  - 4:     **return** 0
  - 5: Use  $A_{jth-rev}^{\mathcal{I} \cup \{u\}}$  to get the number  $\hat{j}$ , such that  $b$  is the  $\hat{j}$ th element in the enumeration of  $S_{\mathcal{I} \cup \{u\}}$ .
  - 6: **return**  $\hat{j}$
- 

If such an element does not exist,  $\left| S_i \cap \tilde{S}_u \right| = 0$ . Since  $A_{enum}^{\mathcal{I} \cup \{u\}}$  enumerates the elements in lexicographical order, the number we obtain in  $A_{jth-rev}^{\mathcal{I} \cup \{u\}}$  is equal to the cardinality of  $S_{\mathcal{I}} \cap \tilde{S}_u$  (since we enumerate all elements in  $c \in \{a_1, \dots, a_{P_u(j)}\}$  with  $c \in S_{\mathcal{I}}$  until  $b$  will be enumerated).

Algorithm 17 takes time  $O(t_j + t_l + t_{jr})$  to compute cardinalities of sets of the form  $\bigcap_{i \in \mathcal{I}} S_i \cap \tilde{S}_u$ . The number of sets of such a form in the formula of  $P_u(j)$  (see Equation 9.8 and 9.9) is at most  $(k-1)2^{k-1}$ . All in all, the algorithm takes time  $O((k-1)2^{k-1}(t_j + t_l + t_{jr} + t_c))$ .

To execute Algorithm 16, we have to compute  $P_k(j)$  for all  $k \in \{\ell, \dots, n\}$ . Therefore, it takes time at most  $O(m(m-1)2^{m-1}(t_j + t_l + t_{jr} + t_c))$  to compute all these values.

In summary, Algorithm 16 takes time

$$\begin{aligned}
& O(m(m-1)2^{m-1}(t_j + t_l + t_{jr} + t_c) + \underbrace{m2^m t_c}_{\text{line 4}} + \underbrace{t_j}_{\text{line 5}} + \underbrace{m(t_l + t_j)}_{\text{while-loop}}) \\
& = O(m^2 2^m (t_j + t_l + t_{jr} + t_c) + m(t_l + t_j))
\end{aligned}$$

This concludes the proof of Lemma 9.3.

In the remainder of this section we will prove Theorem 3.12 where we apply Theorem 3.9 on Lemma 9.1 and 9.3.

To use Theorem 3.9 on Lemma 9.1, we will show that for a strongly exhaustively  $q$ -hierarchical query  $\bigcup_{i=1}^m Q_i$  it holds that for every  $\emptyset \neq \mathcal{I} \subseteq [m]$  the query  $Q_{\mathcal{I}}$  that expresses that

$$Q_{\mathcal{I}}(D) = \bigcap_{i \in \mathcal{I}} Q_i(D) \text{ for all } \sigma\text{-dbs } D$$

is also  $q$ -hierarchical. This fact follows from the following claim

**Claim 9.6.** *Let the query*

$$Q := \{(x_1, \dots, x_k, b_1, \dots, b_\ell) : \exists y_1 \dots \exists y_\ell R_1 \bar{z}_1 \wedge \dots \wedge R_m \bar{z}_m\}$$

*be a  $q$ -hierarchical query. Then, the query*

$$Q' := \{(x_{i_1}, \dots, x_{i_{k'}}, b_{u_1}, \dots, b_{u_{\ell'}}) : \exists y_{j_1} \dots \exists y_{j_{\ell'}} R_1 \bar{z}_1 \wedge \dots \wedge R_{m-1} \bar{z}_{m-1}\}$$

## 9. Answering unions of conjunctive queries under updates

where

- $\{x_{i_1}, \dots, x_{i_{k'}}\} = \{x_1, \dots, x_k\} \cap \bigcup_{M \in [m-1]} \{z : z \text{ is in } \bar{z}_M\}$  and
- $\{y_{j_1}, \dots, y_{j_{\ell'}}\} = \{y_1, \dots, y_\ell\} \cap \bigcup_{M \in [m-1]} \{z : z \text{ is in } \bar{z}_M\}$  and
- $\{b_{u_1}, \dots, b_{u_{\ell'}}\} = \{b_1, \dots, b_\ell\} \cap \bigcup_{M \in [m-1]} \{z : z \text{ is in } \bar{z}_M\}$  and

let  $T$  be a  $q$ -tree of  $Q$  and  $T'$  be the induced subtree of  $T$  on  $\text{vars}(Q')$ . Then,  $T'$  is a  $q$ -tree for  $Q'$ . In particular,  $Q'$  is  $q$ -hierarchical.

In particular, we obtain the query  $Q_{\mathcal{I}}$  by removing atoms from the query  $Q_{[m]}$ .

Before we prove Claim 9.6 we show how to use the fact to apply Theorem 3.9 on Lemma 9.1 and 9.3 to obtain proof for Theorem 3.12.

*Proof for Theorem 3.12.* We use in parallel for every  $\emptyset \neq \mathcal{I} \subseteq [m]$  the data structures  $D_{\mathcal{I}}$  from Theorem 3.9 for  $Q_{\mathcal{I}}$  and  $D$ .

Whenever the database receives an update, we update the data structures  $D_{\mathcal{I}}$  in time  $\text{poly}(Q) \cdot O(\log(\|D\|))$  for all  $\emptyset \neq \mathcal{I} \subseteq [m]$ . All in all, this can be done in time  $\exp(Q) \cdot O(\log(\|D\|))$ .

From Theorem 3.9 we have for all  $\emptyset \neq \mathcal{I} \subseteq [m]$  corresponding algorithms available (obtained from the assumption of Lemma 9.3). For all  $\emptyset \neq \mathcal{I} \subseteq [m]$  the algorithm  $A_{\text{enum}}^{\mathcal{I}}$  has delay  $t_d = O(\text{poly}(Q))$  (see Theorem 3.9(b)),  $A_{\text{count}}^{\mathcal{I}}$  can be done in time  $t_c = O(1)$  (see Theorem 3.9(a)),  $A_{\text{test}}^{\mathcal{I}}$  can be done in time  $t_t = O(\text{poly}(Q))$  (see Theorem 3.9(c)),  $A_{\text{jth}}^{\mathcal{I}}$  can be done in time  $t_j = \text{poly}(Q) \cdot O(\log(\|D\|))$  (see Theorem 3.9(e)),  $A_{\text{jth-rev}}^{\mathcal{I}}$  can be done in time  $t_{jr} = \text{poly}(Q) \cdot O(\log(\|D\|))$  (see Theorem 3.9(f)),  $A_{\text{lex}}^{\mathcal{I}}$  can be done in time  $t_l = \text{poly}(Q)$  (see Theorem 3.9(d)).

To enumerate the tuples in  $\bigcup_{i=1}^m Q_i(D)$ , we use the enumeration algorithm from Lemma 9.1 with the enumeration algorithms  $A_{\text{enum}}^{\{i\}}$  and testing algorithms  $A_{\text{test}}^{\{i\}}$  for all  $i \in [m]$ . Clearly, the tuples can be enumerated with delay  $\text{poly}(Q)$ .

Furthermore, we can apply Lemma 9.3 to obtain an algorithm that upon input of a number  $j \in \mathbb{N}$  outputs the tuple  $a \in \bigcup_{i=1}^m Q_i$  such that  $a$  will be the  $j$ th tuple in the enumeration in time  $\exp(Q)O(\log(\|D\|))$ .

This concludes the proof of Theorem 3.12.  $\square$

In the remainder of this section we prove Claim 9.6:

*Proof of Claim 9.6.* We have to show that the following conditions holds:

- (1) for all atoms  $\psi$  in  $Q'$  the set  $\text{vars}(\psi)$  forms a directed path in  $T'$  that starts from the root and
- (2) if  $\text{free}(Q') \neq \emptyset$ , then  $\text{free}(Q')$  is a connected subset in  $T'$  that contains the root.

Let  $\psi$  be an arbitrary atom in  $Q'$ . Clearly, since  $\psi$  is an atom in  $Q$  and  $Q$  is  $q$ -hierarchical, the set  $\text{vars}(\psi)$  forms a directed path in  $T$  (and since  $\text{vars}(\psi) \subseteq \text{vars}(Q')$ , in particular, in  $T'$ ) that starts from the root. Therefore, Condition (1) holds.



## 9.2. Reporting the $j$ th solution

To show that Condition (2) holds, we assume for contradiction that there are (at least) two connected components in  $\text{free}(Q')$  in  $T'$ . Let the free variables  $z_1$  and  $z_2$  be arbitrary variables of the connected components. Since  $z_1, z_2 \in V(T') = \text{vars}(Q')$ , there are  $M, M' \in [m - 1]$  such that  $z_1$  is in  $\bar{z}_M$  and  $z_2$  is in  $\bar{z}_{M'}$ . Since  $Q$  is q-hierarchical there is a path in  $T$  (and in  $T'$ ) from the root to  $z_1$  and from the root to  $z_2$  that only use free variables. These paths only contain variables in  $(\{z : z \text{ is in } \bar{z}_M\} \cap \text{free}(Q')) \cup (\{z : z \text{ is in } \bar{z}_{M'}\} \cap \text{free}(Q'))$ . In particular there is a path from the root to  $z_1$  and a path from the root to  $z_2$  and therefore,  $z_1$  and  $z_2$  do not belong to different connected components.  $\square$



**Part II.**

**Answering FO+MOD Queries  
under Updates on Bounded  
Degree Databases**



## 10. Further preliminaries and main results in Part II

This part of the thesis is based on [20]. Here, we consider the task of answering FO+MOD queries under updates on bounded degree databases.

In this section we start with further preliminaries that are necessary for this part. Then the second part's main result is given and at the end the organization of this part is given.

**Gaifman graph.** The *Gaifman graph* of a  $\sigma$ -db  $D$  is the undirected simple graph  $G^D = (V, E)$  with vertex set  $V := \text{adom}(D)$ , where there is an edge between vertices  $u$  and  $v$  whenever  $u \neq v$  and there are  $R \in \sigma$  and  $(a_1, \dots, a_{\text{ar}(R)}) \in R^D$  such that  $u, v \in \{a_1, \dots, a_{\text{ar}(R)}\}$ . A  $\sigma$ -db  $D$  is called *connected* if its Gaifman graph  $G^D$  is connected; the *connected components* of  $D$  are the connected components of  $G^D$ . The *degree* of a database  $D$  is the degree of its Gaifman graph  $G^D$ , i.e., the maximum number of neighbours of a node of  $G^D$ .

Throughout this part of the thesis we fix a number  $d \in \mathbb{N}$  and restrict attention to *d-bounded* databases, i.e., to databases of degree at most  $d$ .

**Further notes for updates.** Note that update commands may change the database's degree. In this part of this thesis, we restrict attention to databases of degree at most  $d$ . Therefore, when applying an insertion command to a  $\sigma$ -db  $D$  of degree  $\leq d$ , the command is carried out only if the resulting database  $D'$  still has degree  $\leq d$ ; otherwise  $D$  remains unchanged and instead of carrying out the insertion command, an error message is returned.

**Queries.** We fix a countably infinite set **var** of *variables*. In this part of the thesis, we consider the extension FO+MOD of first-order logic FO with modulo-counting quantifiers. For a fixed schema  $\sigma$ , the set FO+MOD[ $\sigma$ ] is built from atomic formulas of the form  $x_1 = x_2$  and  $R(x_1, \dots, x_{\text{ar}(R)})$ , for  $R \in \sigma$  and variables  $x_1, x_2, \dots, x_{\text{ar}(R)} \in \mathbf{var}$ , and is closed under Boolean connectives  $\neg, \wedge$ , existential first-order quantifiers  $\exists x$ , and modulo-counting quantifiers  $\exists^{i \bmod m} x$ , for a variable  $x \in \mathbf{var}$  and integers  $i, m \in \mathbb{N}$  with  $m \geq 2$  and  $i < m$ . The intuitive meaning of a formula of the form  $\exists^{i \bmod m} x \psi$  is that the number of witnesses  $x$  that satisfy  $\psi$  is congruent  $i$  modulo  $m$ . Note that FO+MOD is strictly more expressive than first-order logic without counting quantifiers, since it can express that the number of elements in a unary relation is even, which is not possible to express in FO (cf., e.g., [72]). As usual,  $\forall x, \vee, \rightarrow, \leftrightarrow$  will be used as abbreviations when constructing formulas. It will be convenient to add the quantifiers  $\exists^{\geq m} x$ , for  $m \in \mathbb{N}_{\geq 1}$ ; a formula of the form  $\exists^{\geq m} x \psi$  expresses that the number of witnesses  $x$  which satisfy  $\psi$  is  $\geq m$ . Though these quantifiers allow more

## 10. Further preliminaries and main results in Part II

succinct definitions, we will treat them as syntactic sugar since they do not increase the expressive power of FO+MOD.

The *quantifier rank*  $qr(\varphi)$  of a FO+MOD-formula  $\varphi$  is the maximum nesting depth of quantifiers that occur in  $\varphi$ . By  $\text{free}(\varphi)$  we denote the set of all *free variables* of  $\varphi$ , i.e., all variables  $x$  that have at least one occurrence in  $\varphi$  that is not within a quantifier of the form  $\exists x$ ,  $\exists^{\geq m} x$ , or  $\exists^{i \bmod m} x$ . A *sentence* is a formula  $\varphi$  with  $\text{free}(\varphi) = \emptyset$ .

An *assignment* for  $\varphi$  is a partial mapping  $\alpha$  from **var** to **dom**, where  $\text{free}(\varphi) \subseteq \text{dom}(\alpha)$ . We write  $(D, \alpha) \models \varphi$  to indicate that  $\varphi$  is satisfied when evaluated in  $D$  with respect to *active domain semantics* while interpreting every free occurrence of a variable  $x$  with the constant  $\alpha(x)$ . Recall from [2] that “active domain semantics” means that quantifiers are evaluated with respect to the database’s active domain. In particular,  $(D, \alpha) \models \exists x \psi$  if and only if there exists an  $a \in \text{adom}(D)$  such that  $(D, \alpha_x^a) \models \psi$ , where  $\alpha_x^a$  is the assignment  $\alpha'$  with  $\alpha'(x) = a$  and  $\alpha'(y) = \alpha(y)$  for all  $y \in \text{dom}(\alpha) \setminus \{x\}$ . Accordingly,  $(D, \alpha) \models \exists^{\geq m} x \psi$  if and only if  $|\{a \in \text{adom}(D) : (D, \alpha_x^a) \models \psi\}| \geq m$ , and  $(D, \alpha) \models \exists^{i \bmod m} x \psi$  if and only if  $|\{a \in \text{adom}(D) : (D, \alpha_x^a) \models \psi\}| \equiv i \bmod m$ .

A  $k$ -ary FO+MOD query of schema  $\sigma$  is of the form  $Q(x_1, \dots, x_k)$  where  $k \in \mathbb{N}$ ,  $Q \in \text{FO+MOD}[\sigma]$ , and  $\text{free}(Q) \subseteq \{x_1, \dots, x_k\}$ . We will often assume that the tuple  $(x_1, \dots, x_k)$  is clear from the context and simply write  $Q$  instead of  $Q(x_1, \dots, x_k)$  and  $(D, (a_1, \dots, a_k)) \models Q$  instead of  $(D, \frac{a_1, \dots, a_k}{x_1, \dots, x_k}) \models Q$ , where  $\frac{a_1, \dots, a_k}{x_1, \dots, x_k}$  denotes the assignment  $\alpha$  with  $\alpha(x_i) = a_i$  for all  $i \in [k]$ . When evaluated in a  $\sigma$ -db  $D$ , the  $k$ -ary query  $Q(x_1, \dots, x_k)$  yields the  $k$ -ary relation

$$Q(D) := \{(a_1, \dots, a_k) \in \text{adom}(D)^k : (D, \frac{a_1, \dots, a_k}{x_1, \dots, x_k}) \models Q\}.$$

*Boolean queries* are  $k$ -ary queries with  $k = 0$ . As usual, for Boolean queries we will write  $Q(D) = \text{no}$  instead of  $Q(D) = \emptyset$ , and  $Q(D) = \text{yes}$  instead of  $Q(D) \neq \emptyset$ ; and we write  $D \models Q$  to indicate that  $(D, \alpha) \models Q$  for any assignment  $\alpha$ .

**Sizes and Cardinalities of Queries.** The *size*  $\|\sigma\|$  of a schema  $\sigma$  is the sum of the arities of its relation symbols. The size  $\|Q\|$  of an FO+MOD query  $Q$  of schema  $\sigma$  is the length of  $Q$  when viewed as a word over the alphabet  $\sigma \cup \mathbf{var} \cup \mathbb{N} \cup \{=, \wedge, \neg, \exists, \text{mod}, \geq, (, )\} \cup \{, \}$ . For a  $k$ -ary query  $Q(x_1, \dots, x_k)$  and a  $\sigma$ -db  $D$ , the *cardinality of the query result* is the number  $|Q(D)|$  of tuples in  $Q(D)$ . The *cardinality*  $|D|$  of a  $\sigma$ -db  $D$  is defined as the number of tuples stored in  $D$ , i.e.,  $|D| := \sum_{R \in \sigma} |R^D|$ . The *size*  $\|D\|$  of  $D$  is defined as  $\|\sigma\| + |\text{adom}(D)| + \sum_{R \in \sigma} \text{ar}(R) \cdot |R^D|$  and corresponds to the size of a reasonable encoding of  $D$ . Throughout this part of the thesis we let  $f(Q, d)$  stand for a function of the form

$$f(Q, d) = 2^{d^{2^{O(\|Q\|)}}}. \quad (10.1)$$

For our purposes it will be convenient to assume that **dom** =  $\mathbb{N}_{\geq 1}$ .

Our algorithms in this part will take as input a  $k$ -ary FO+MOD-query  $Q(x_1, \dots, x_k)$ , a parameter  $d$ , and an initial  $\sigma$ -db  $D_0$  of degree  $\leq d$ .

**Hanf Normal Form for FO+MOD.** Our algorithms for evaluating FO+MOD queries rely on a decomposition of FO+MOD queries into *Hanf normal form*. To describe this normal form, we need some more notation.

Two formulas  $\varphi$  and  $\psi$  of schema  $\sigma$  are called *d-equivalent* (in symbols:  $\varphi \equiv_d \psi$ ) if for all  $\sigma$ -dbs  $D$  of degree  $\leq d$  and all assignments  $\alpha$  for  $\varphi$  and  $\psi$  in  $D$  we have  $(D, \alpha) \models \varphi \iff (D, \alpha) \models \psi$ .

For a  $\sigma$ -db  $D$  and a set  $A \subseteq \text{adom}(D)$  we write  $D[A]$  to denote the restriction of  $D$  to the domain  $A$ , i.e.,  $R^{D[A]} = \{\bar{a} \in R^D : \bar{a} \in A^{\text{ar}(R)}\}$ , for all  $R \in \sigma$ . For two  $\sigma$ -dbs  $D$  and  $D'$ , an *isomorphism*  $\pi : D \rightarrow D'$  is a bijection from  $\text{adom}(D)$  to  $\text{adom}(D')$  with  $(b_1, \dots, b_r) \in R^D \iff (\pi(b_1), \dots, \pi(b_r)) \in R^{D'}$  for all  $R \in \sigma$ , for  $r := \text{ar}(R)$ , and for all  $b_1, \dots, b_r \in \text{adom}(D)$ . For two  $k$ -tuples  $\bar{a} = (a_1, \dots, a_k)$  and  $\bar{a}' = (a'_1, \dots, a'_k)$  of elements in  $\text{adom}(D)$  and  $\text{adom}(D')$ , resp., we write  $(D, \bar{a}) \cong (D', \bar{a}')$  to indicate that there is an isomorphism  $\pi$  from  $D$  to  $D'$  that maps  $a_i$  to  $a'_i$  for all  $i \in [k]$ .

The *distance*  $\text{dist}^D(a, b)$  between two elements  $a, b \in \text{adom}(D)$  is the minimal length (i.e., the number of edges) of a path from  $a$  to  $b$  in  $D$ 's Gaifman graph  $G^D$  (if no such path exists, we let  $\text{dist}^D(a, b) = \infty$ ; note that  $\text{dist}^D(a, a) = 0$ ). For  $r \geq 0$  and  $a \in \text{adom}(D)$ , the *r-ball* around  $a$  in  $D$  is the set

$$N_r^D(a) := \{b \in \text{adom}(D) : \text{dist}^D(a, b) \leq r\}.$$

For a  $\sigma$ -db  $D$  and a tuple  $\bar{a} = (a_1, \dots, a_k)$  we let  $N_r^D(\bar{a}) := \bigcup_{i \in [k]} N_r^D(a_i)$ . The *r-neighbourhood* around  $\bar{a}$  in  $D$  is defined as the  $\sigma$ -db  $\mathcal{N}_r^D(\bar{a}) := D[N_r^D(\bar{a})]$ . For  $r \geq 0$  and  $k \geq 1$ , a *type*  $\tau$  (over  $\sigma$ ) with  $k$  centres and radius  $r$  (for short: *r-type with k centres*) is of the form  $(T, \bar{t})$ , where  $T$  is a  $\sigma$ -db,  $\bar{t} \in \text{adom}(T)^k$ , and  $\text{adom}(T) = N_r^T(\bar{t})$ . The elements in  $\bar{t}$  are called the *centres* of  $\tau$ . For a tuple  $\bar{a} \in \text{adom}(D)^k$ , the *r-type of  $\bar{a}$  in  $D$*  is defined as the *r-type with k centres*  $(\mathcal{N}_r^D(\bar{a}), \bar{a})$ .

For a given *r-type with k centres*  $\tau = (T, \bar{t})$  it is straightforward to construct a first-order formula  $\text{sph}_\tau(\bar{x})$  (depending on  $r$  and  $\tau$ ) with  $k$  free variables  $\bar{x} = (x_1, \dots, x_k)$  which expresses that the *r-type of  $\bar{x}$*  is isomorphic to  $\tau$ , i.e., for every  $\sigma$ -db  $D$  and all  $\bar{a} = (a_1, \dots, a_k) \in \text{adom}(D)^k$  we have

$$(D, \bar{a}) \models \text{sph}_\tau(\bar{x}) \iff (\mathcal{N}_r^D(\bar{a}), \bar{a}) \cong (T, \bar{t}).$$

The formula  $\text{sph}_\tau(\bar{x})$  is called a *sphere-formula* (over  $\sigma$  and  $\bar{x}$ ); the numbers  $r$  and  $k$  are called *locality radius* and *arity*, resp., of the sphere-formula.

A *Hanf-sentence* (over  $\sigma$ ) is a sentence of the form

$$\exists^{\geq m} x \text{sph}_\tau(x) \quad \text{or} \quad \exists^{i \bmod m} x \text{sph}_\tau(x),$$

where  $\tau$  is an *r-type* (over  $\sigma$ ) with 1 centre, for some  $r \geq 0$ . The number  $r$  is called *locality radius* of the Hanf-sentence. A formula in *Hanf normal form* (over  $\sigma$ ) is a Boolean combination<sup>1</sup> of sphere-formulas and Hanf-sentences (over  $\sigma$ ). The *locality radius* of a formula  $\psi$  in Hanf normal form is the maximum of the locality radii of the Hanf-sentences and the sphere-formulas that occur in  $\psi$ . The formula is *d-bounded* if all types  $\tau$  that occur in sphere-formulas or Hanf-sentences of  $\psi$  are *d-bounded*, i.e.,  $T$  is of degree  $\leq d$ , where  $\tau = (T, \bar{t})$ . Our query evaluation algorithms for FO+MOD rely on the following result by Heimberg, Kuske, and Schweikardt [55].

<sup>1</sup>Throughout this thesis, whenever we speak of *Boolean combinations* we mean *finite* Boolean combinations.

## 10. Further preliminaries and main results in Part II

**Theorem 10.1** ([55]). *There is an algorithm which receives as input a degree bound  $d \in \mathbb{N}$  and an  $\text{FO}+\text{MOD}[\sigma]$ -formula  $\varphi$ , and constructs a  $d$ -equivalent formula  $\psi$  in Hanf normal form (over  $\sigma$ ) with the same free variables as  $\varphi$ . For any  $d \geq 2$ , the formula  $\psi$  is  $d$ -bounded and has locality radius  $\leq 4^{\text{qr}(\varphi)}$ , and the algorithm's runtime is  $2^{d^{2^{O(\|\varphi\|+\|\sigma\|)}}}$ .*

The first step of all our query evaluation algorithms is to use Theorem 10.1 to transform a given query  $Q(\bar{x})$  into a  $d$ -equivalent query  $\psi(\bar{x})$  in Hanf normal form. The following lemma summarises standard facts, that we will apply at several places throughout Part II to evaluate the sphere-formulas that occur in  $\psi$ .

**Lemma 10.2.** *Let  $d \geq 2$  and let  $D$  be a  $\sigma$ -db of degree  $\leq d$ . Let  $r \geq 0$ ,  $k \geq 1$ , and  $\bar{a} = (a_1, \dots, a_k) \in \text{adom}(D)^k$ .*

- (a)  $|N_r^D(\bar{a})| \leq k \sum_{i=0}^r d^i \leq kd^{r+1}$ .
- (b) *Given  $D$  and  $\bar{a}$ , the  $r$ -neighbourhood  $N_r^D(\bar{a})$  can be computed in time  $(kd^{r+1})^{O(\|\sigma\|)}$ .*
- (c)  $N_r^D(a_1, a_2)$  is connected if and only if  $\text{dist}^D(a_1, a_2) \leq 2r+1$ .
- (d) *If  $N_r^D(\bar{a})$  is connected, then  $N_r^D(\bar{a}) \subseteq N_{r+(k-1)(2r+1)}^D(a_i)$ , for all  $i \in [k]$ .*
- (e) *Let  $D'$  be a  $\sigma$ -db of degree  $\leq d$  and let  $\bar{b} = (b_1, \dots, b_k) \in \text{adom}(D')^k$ . It can be tested in time  $(kd^{r+1})^{O(\|\sigma\|+kd^{r+1})} \leq 2^{O(\|\sigma\|k^2d^{2r+2})}$  whether  $(N_r^D(\bar{a}), \bar{a}) \cong (N_r^{D'}(\bar{b}), \bar{b})$ .*

*Proof.* Parts (a)–(d) are straightforward. Concerning Part (e), a brute-force approach is to loop through all mappings from  $N_r^D(\bar{a})$  to  $N_r^{D'}(\bar{b})$  that map  $a_i$  to  $b_i$  for every  $i \in [k]$  and check whether this mapping is an isomorphism. Each such check can be accomplished in time  $n^{O(\|\sigma\|)}$  for  $n := kd^{r+1}$ , and the number of mappings that have to be checked is  $\leq n^n$ . Thus, the isomorphism test is accomplished in time  $n^{O(n+\|\sigma\|)} = (kd^{r+1})^{O(\|\sigma\|+kd^{r+1})}$ .  $\square$

The time bound stated in part (e) of Lemma 10.2 is obtained by a brute-force approach. When using Luks' polynomial time isomorphism test for bounded degree graphs [75], the time bound of Lemma 10.2(e) can be improved to  $(kd^{r+1})^{\text{poly}(d\|\sigma\|)}$ . However, the asymptotic overall runtime of our algorithms for evaluating  $\text{FO}+\text{MOD}$ -queries won't improve when using Luks' algorithm instead of the brute-force isomorphism test of Lemma 10.2(e).

**Main Results.** The aim of this part is to show the following:

On input of an initial database  $D_0$ , and a degree bound  $d$ , and a  $\text{FO}+\text{MOD}$ -query  $Q$  we can construct in a linear time  $f(\|Q\|, d)\|D\|$  preprocessing phase a data structure that can be updated in time  $f(\|Q\|, d)$  and allows to

- immediately answer  $Q$  on  $D$  if  $Q$  is a Boolean query (Theorem 11.1),



- test for a given tuple whether it belongs to the result set in time  $O(k^2)$  (Theorem 12.5),
- immediately output the number of tuples in the result set  $Q(D)$  (Theorem 12.11),
- enumerate the tuples in the result set  $Q(D)$  with delay  $O(k^2)$  (Theorem 12.14), and
- enumerate the tuples in  $Q(D^-) \setminus Q(D^+)$  and  $Q(D^+) \setminus Q(D^-)$  with delay  $O(k^2)$ , where  $D^-$  and  $D^+$  denote the database before and after performing the update operation, respectively (Theorem 12.24).

**Organization.** After some basic definitions in this chapter we obtain a dynamic algorithm for Boolean FO+MOD-queries in Chapter 11. After some preparations for non-Boolean queries in Chapter 12.1, we present the algorithm for testing in Chapter 12.2. In Chapter 12.3 we reduce the task of counting and enumerating FO+MOD-queries in the dynamic setting to the problem of counting and enumerating independent sets in graphs of bounded degree. We use this reduction to provide efficient dynamic counting and enumeration algorithms in Chapter 12.4 and 12.5, respectively. In Chapter 12.6 we generalise this to be able to efficiently enumerate particular subsets of the query result, and we use this in Chapter 12.7 to obtain an efficient dynamic algorithm for enumerating the difference between the old and the new query result. We conclude in Chapter 13.



# 11. Answering Boolean FO+MOD Queries Under Updates

In [44], Frick and Grohe showed that in the static setting (i.e., without database updates), Boolean FO-queries  $Q$  can be answered on databases  $D$  of degree  $\leq d$  in time  $f(Q, d) \cdot \|D\|$ . Our first main theorem extends their result to FO+MOD-queries and the dynamic setting.

**Theorem 11.1.** *There is a dynamic algorithm that receives a schema  $\sigma$ , a degree bound  $d \geq 2$ , a Boolean FO+MOD $[\sigma]$ -query  $Q$  and a  $\sigma$ -db  $D_0$  of degree  $\leq d$ , and computes within  $t_p = f(Q, d) \cdot \|D_0\|$  preprocessing time a data structure that can be updated in time  $t_u = f(Q, d)$  and allows to return the query result  $Q(D)$  with answer time  $t_{ans} = O(1)$ .*

*If  $Q$  is a  $d$ -bounded Hanf-sentence of locality radius  $r$ , then  $f(Q, d)$  improves to  $f(Q, d) = 2^{O(\|\sigma\|d^{2r+2})}$ , and the initialisation time is  $t_i = O(\|Q\|)$ .*

The proof will be an easy consequence of Theorem 10.1 and the following lemma.

**Lemma 11.2.** *There is a dynamic algorithm that receives a schema  $\sigma$ , a number  $s \in \mathbb{N}_{\geq 1}$ , a list  $(r_j)_{j \in [s]}$  of non-negative integers, a list  $(\rho_j)_{j \in [s]}$  where each  $\rho_j$  is a  $r_j$ -type with 1 centre (over  $\sigma$ ), a degree bound  $d \geq 2$  and a  $\sigma$ -db  $D$  of degree  $\leq d$ . The algorithm computes within  $t_i = O(s)$  initialisation time a data structure that can be updated in time  $t_u = \sum_{j=1}^s (d^{r_j+1})^{O(\|\sigma\|+d^{r_j+1})}$ . Upon input of a number  $j \in [s]$  the algorithm returns within time  $O(1)$  the number  $\left| \left\{ a \in \text{adom}(D) : (\mathcal{N}_{r_j}^D(a), a) \cong \rho_j \right\} \right|$ .*

*In particular, the update time  $t_u$  is at most  $s \cdot 2^{O(\|\sigma\| \cdot d^{2r+2})}$ , for  $r := \max_{j \in [s]} r_j$ .*

*Proof.* For each  $j \in [s]$  our data structure will store the number  $A[j]$  of all elements  $a \in \text{adom}(D)$  whose  $r_j$ -type is isomorphic to  $\rho_j$ , i.e.,  $(\mathcal{N}_{r_j}^D(a), a) \cong \rho_j$ . The initialisation for the empty database  $D_\emptyset$  lets  $A[j] = 0$  for all  $j \in [s]$ .

To update our data structure upon a command  $\text{update } R(a_1, \dots, a_k)$ , for  $k = \text{ar}(R)$  and  $\text{update} \in \{\text{insert}, \text{delete}\}$ , we proceed as follows. The idea is to remove from the data structure the information on all the database elements whose  $r_j$ -neighbourhood (for some  $j \in [s]$ ) is affected by the update, and then to recompute the information concerning all these elements on the updated database.

Let  $D^-$  be the database before the update is received and let  $D^+$  be the database after the update has been performed. We consider each  $j \in [s]$ . All elements whose  $r_j$ -neighbourhood might have changed, belong to the set  $U_j := N_{r_j}^{D'}(\bar{a})$ , where  $D' := D^+$  if the update command is  $\text{insert } R(\bar{a})$ , and  $D' := D^-$  if the update command is  $\text{delete } R(\bar{a})$ .

## 11. Answering Boolean FO+MOD Queries Under Updates

To remove the old information from  $\mathbf{A}[j]$ , we compute for each  $a \in U_j$  the neighbourhood  $T_a := \mathcal{N}_{r_j}^{D^-}(a)$ , check whether  $(T_a, a) \cong \rho_j$ , and if so, decrement the value  $\mathbf{A}[j]$ .

To recompute the new information for  $\mathbf{A}[j]$ , we compute for all  $a \in U_j$  the neighbourhood  $T'_a := \mathcal{N}_{r_j}^{D^+}(a)$ , check whether  $(T'_a, a) \cong \rho_j$ , and if so, increment the value  $\mathbf{A}[j]$ .

Using Lemma 10.2 we obtain for each  $j \in [s]$  that  $|U_j| \leq kd^{r_j+1}$ . For each  $a \in U_j$ , the neighbourhoods  $T_a$  and  $T'_a$  can be computed in time  $(d^{r_j+1})^{O(\|\sigma\|)}$ , and testing for isomorphism with  $\rho_j$  can be done in time  $(d^{r_j+1})^{O(\|\sigma\|+d^{r_j+1})}$ . Thus, the update of  $\mathbf{A}[j]$  is done in time  $k \cdot (d^{r_j+1})^{O(\|\sigma\|+d^{r_j+1})}$ . Recall that  $k = \text{ar}(R) \leq \|\sigma\|$ . Hence, the entire update time is  $t_u = \sum_{j=1}^s (d^{r_j+1})^{O(\|\sigma\|+d^{r_j+1})}$ . Finally, note that

$$(d^{r+1})^{O(\|\sigma\|+d^{r+1})} \leq 2^{d^{r+1} \cdot O(\|\sigma\|+d^{r+1})} \leq 2^{O(\|\sigma\| \cdot d^{2r+2})}.$$

This completes the proof of Lemma 11.2.  $\square$

*Proof of Theorem 11.1.*

W.l.o.g. we assume that all the symbols of  $\sigma$  occur in  $Q$  (otherwise, we remove from  $\sigma$  all symbols that do not occur in  $Q$ ). In the preprocessing routine, we first use Theorem 10.1 to transform  $Q$  into a  $d$ -equivalent sentence  $\psi$  in Hanf normal form; this takes time  $f(Q, d)$ . The sentence  $\psi$  is a Boolean combination of  $d$ -bounded Hanf-sentences (over  $\sigma$ ) of locality radius at most  $r := 4^{\text{ar}(Q)}$ . Let  $\rho_1, \dots, \rho_s$  be the list of all types that occur in  $\psi$ . Thus, every Hanf-sentence in  $\psi$  is of the form  $\exists^{\geq k} x \text{ sph}_{\rho_j}(x)$  or  $\exists^{i \bmod m} x \text{ sph}_{\rho_j}(x)$  for some  $j \in [s]$  and  $k, i, m \in \mathbb{N}$  with  $k \geq 1$ ,  $m \geq 2$ , and  $i < m$ . For each  $j \in [s]$  let  $r_j$  be the radius of  $\text{sph}_{\rho_j}(x)$ . Thus,  $\rho_j$  is an  $r_j$ -type with 1 centre (over  $\sigma$ ).

We use the dynamic data structure provided by Lemma 11.2, and in addition, we also store a Boolean value  $\mathbf{Ans}$  where  $\mathbf{Ans} = Q(D)$  is the answer of the Boolean query  $Q$  on the current database  $D$ . This way, the query can be answered in time  $O(1)$  by simply outputting  $\mathbf{Ans}$ .

The initialisation for the empty database  $D_\emptyset$  computes  $\mathbf{Ans}$  as follows. Every Hanf-sentence of the form  $\exists^{\geq k} x \text{ sph}_{\rho_j}(x)$  in  $\psi$  is replaced by the Boolean constant **false**. Every Hanf-sentence of the form  $\exists^{i \bmod m} x \text{ sph}_{\rho_j}(x)$  is replaced by **true** if  $i = 0$  and by **false** otherwise. The resulting formula, a Boolean combination of the Boolean constants **true** and **false**, then is evaluated, and we let  $\mathbf{Ans}$  be the obtained result. The entire initialisation takes time at most  $t_i = \|\psi\| = f(Q, d) = 2^{d^{2O(\|Q\|)}}$ .

To update our data structure upon a command  $\text{update } R(a_1, \dots, a_k)$ , we first perform the update routine of the data structure provided by Lemma 11.2. Afterwards, we recompute the query answer  $\mathbf{Ans}$  as follows. Every Hanf-sentence of the form  $\exists^{\geq k} x \text{ sph}_{\rho_j}(x)$  in  $\psi$  is replaced by the Boolean constant **true** if

$$\left| \left\{ a \in \text{adom}(D) : (\mathcal{N}_{r_j}^D(a), a) \cong \rho_j \right\} \right| \geq k,$$

and by the Boolean constant **false** otherwise. Every Hanf-sentence of the form  $\exists i \bmod m x \text{ sph}_{\rho_j}(x)$  is replaced by **true** if

$$\left| \left\{ a \in \text{adom}(D) : (\mathcal{N}_{r_j}^D(a), a) \cong \rho_j \right\} \right| \equiv i \bmod m,$$

and by **false** otherwise. The resulting formula, a Boolean combination of the Boolean constants **true** and **false**, then is evaluated, and we let **Ans** be the obtained result. Thus, recomputing **Ans** takes time  $\text{poly}(\|\psi\|)$ .

Noting that  $r_j \leq 4^{\text{qr}(Q)} \leq 2^{O(\|Q\|)}$  and  $s \leq \|\psi\|$ , we obtain that the entire update time is

$$t_u \leq \sum_{j=1}^s (d^{r_j+1})^{O(\|\sigma\|+d^{r_j+1})} + \text{poly}(\|\psi\|) \leq 2^{d^{2^{O(\|Q\|)}}} = f(Q, d).$$

This completes the proof of Theorem 11.1.  $\square$

In [44], Frick and Grohe obtained a matching lower bound for answering Boolean FO-queries of schema  $\sigma = \{E\}$  on databases of degree at most  $d := 3$  in the static setting. They used the (reasonable) complexity theoretic assumption  $\text{FPT} \neq \text{AW}[*]$  and showed that if this assumption is correct, then there is no algorithm that answers Boolean FO-queries  $Q$  on  $\sigma$ -dbs  $D$  of degree  $\leq 3$  in time  $2^{2^{2^{O(\|Q\|)}}} \cdot \text{poly}(\|D\|)$  in the static setting (see Theorem 2 in [44]). As a consequence, the same lower bound holds in the dynamic setting and shows that in Theorem 11.1, the 3-fold exponential dependency on the query size  $\|Q\|$  cannot be substantially lowered (unless  $\text{FPT} = \text{AW}[*]$ ):

**Corollary 11.3.** *Let  $\sigma := \{E\}$  and let  $d := 3$ . If  $\text{FPT} \neq \text{AW}[*]$ , then there is no dynamic algorithm that receives a Boolean FO $[\sigma]$ -query  $Q$  and a  $\sigma$ -db  $D_0$ , and computes within  $t_p \leq f(Q) \cdot \text{poly}(\|D_0\|)$  preprocessing time a data structure that can be updated in time  $t_u \leq f(Q)$  and allows to return the query result  $Q(D)$  with answer time  $t_{\text{ans}} \leq f(Q)$ , for a function  $f$  with  $f(Q) = 2^{2^{2^{O(\|Q\|)}}}$ .*



## 12. Answering non-Boolean FO+MOD Queries Under Updates

### 12.1. Technical Lemmas on Types and Spheres Useful for Handling Non-Boolean Queries

For our algorithms for evaluating non-Boolean queries it will be convenient to work with a fixed list of representatives of  $d$ -bounded  $r$ -types, provided by the following lemma.

**Lemma 12.1.** *There is an algorithm which upon input of a schema  $\sigma$ , a degree bound  $d \geq 2$ , a radius  $r \geq 0$  and a number  $k \geq 1$ , computes a list  $\mathcal{L}_r^{\sigma,d}(k) = \tau_1, \dots, \tau_\ell$  (for a suitable  $\ell \geq 1$ ) of  $d$ -bounded  $r$ -types with  $k$  centres (over  $\sigma$ ), such that for every  $d$ -bounded  $r$ -type  $\tau$  with  $k$  centres (over  $\sigma$ ) there is exactly one  $i \in [\ell]$  such that  $\tau \cong \tau_i$ . The algorithm's runtime is  $2^{(kd^{r+1})^{O(\|\sigma\|)}}$ . Furthermore, upon input of a  $d$ -bounded  $r$ -type  $\tau$  with  $k$  centres (over  $\sigma$ ), the particular  $i \in [\ell]$  with  $\tau \cong \tau_i$  can be computed in time  $2^{(kd^{r+1})^{O(\|\sigma\|)}}$ .*

Taking into account the statements of Lemma 10.2 (in particular, the time bound provided by Lemma 10.2(e)), the proof of Lemma 12.1 is straightforward. Throughout the remainder of this part,  $\mathcal{L}_r^{\sigma,d}(k)$  will always denote the list provided by Lemma 12.1. The following lemma will be useful for evaluating Boolean combinations of sphere-formulas.

**Lemma 12.2.** *Let  $\sigma$  be a schema, let  $r \geq 0$ ,  $k \geq 1$ ,  $d \geq 2$  and let  $\mathcal{L}_r^{\sigma,d}(k) = \tau_1, \dots, \tau_\ell$ . Let  $\bar{x} = (x_1, \dots, x_k)$  be a list of  $k$  pairwise distinct variables. For every Boolean combination  $\psi(\bar{x})$  of  $d$ -bounded sphere-formulas of radius at most  $r$  (over  $\sigma$ ), there is an  $I \subseteq [\ell]$  such that  $\psi(\bar{x}) \equiv_d \bigvee_{i \in I} \text{sph}_{\tau_i}(\bar{x})$ .*

*Furthermore, given  $\psi(\bar{x})$ , the set  $I$  can be computed in time  $\text{poly}(\|\psi\|) \cdot 2^{(kd^{r+1})^{O(\|\sigma\|)}}$ .*

*Proof.* As a first step, we consider each sphere-formula  $\zeta$  that occurs in  $\psi$  and replace it by a  $d$ -equivalent disjunction of sphere-formulas  $\text{sph}_{\tau_j}(\bar{x})$  with  $\tau_j$  in  $\mathcal{L}_r^{\sigma,d}(k)$ : if  $\zeta$  has arity  $k' \leq k$  and radius  $r' \leq r$  and is of the form  $\text{sph}_\rho(\bar{x}')$  with  $\bar{x}' = (x_{\nu_1}, \dots, x_{\nu_{k'}})$  for  $1 \leq \nu_1 < \dots < \nu_{k'} \leq k$  and  $\rho = (S, \bar{s})$  with  $\bar{s} = (s_1, \dots, s_{k'})$ , then we replace  $\zeta$  by the formula  $\zeta' := \bigvee_{j \in J} \text{sph}_{\tau_j}(\bar{x})$ , where  $J$  consists of all those  $j \in [\ell]$  where for  $(T, \bar{t}) = \tau_j$  with  $\bar{t} = (t_1, \dots, t_k)$  and for  $\bar{t}' := (t_{\nu_1}, \dots, t_{\nu_{k'}})$  we have  $(S, \bar{s}) \cong (T[N_{r'}^T(\bar{t}'), \bar{t}'])$ . It is straightforward to see that  $\zeta'$  and  $\zeta$  are  $d$ -equivalent.

## 12. Answering non-Boolean FO+MOD Queries Under Updates

Let  $\psi_1$  be the formula obtained from  $\psi$  by replacing each  $\zeta$  by  $\zeta'$ . By the Lemmas 12.1 and 10.2,  $\psi_1$  can be constructed in time  $O(\|\psi\| \cdot 2^{(kd^{r+1})^{O(\|\sigma\|)}})$ . Note that  $\psi_1$  is a Boolean combination of formulas  $\text{sph}_{\tau_j}(\bar{x})$  for  $j \in [\ell]$ .

In the second step, we repeatedly use de Morgan's law to push all  $\neg$ -symbols in  $\psi_1$  directly in front of sphere-formulas. Afterwards, we replace every subformula of the form  $\neg \text{sph}_{\tau_j}(\bar{x})$  by the  $d$ -equivalent formula  $\bigvee_{i \in [\ell] \setminus \{j\}} \text{sph}_{\tau_i}(\bar{x})$ . Let  $\psi_2$  be the formula obtained from  $\psi_1$  by these transformations. Constructing  $\psi_2$  from  $\psi_1$  takes time at most  $O(\|\psi_1\|) \cdot 2^{(kd^{r+1})^{O(\|\sigma\|)}} = O(\|\psi\| \cdot 2^{(kd^{r+1})^{O(\|\sigma\|)}})$ .

In the third step, we eliminate all the  $\wedge$ -symbols in  $\psi_2$ . By the definition of the sphere-formulas  $\tau_1, \dots, \tau_\ell$  we have

$$\text{sph}_{\tau_i}(\bar{x}) \wedge \text{sph}_{\tau_{i'}}(\bar{x}) \equiv_d \begin{cases} \text{sph}_{\tau_i}(\bar{x}), & \text{if } i = i' \\ \perp, & \text{if } i \neq i' \end{cases} \quad (12.1)$$

where  $\perp$  is an unsatisfiable formula. Thus, by the distributive law we obtain for all  $m \geq 1$  and all  $I_1, \dots, I_m \subseteq [\ell]$  that

$$\begin{aligned} \bigwedge_{j \in [m]} \left( \bigvee_{i \in I_j} \text{sph}_{\tau_i}(\bar{x}) \right) &\equiv_d \bigvee_{i_1 \in I_1} \dots \bigvee_{i_m \in I_m} \left( \text{sph}_{\tau_{i_1}}(\bar{x}) \wedge \dots \wedge \text{sph}_{\tau_{i_m}}(\bar{x}) \right) \\ &\equiv_d \bigvee_{i \in I} \text{sph}_{\tau_i}(\bar{x}) \end{aligned}$$

for  $I := I_1 \cap \dots \cap I_m$ . We repeatedly use this equivalence during a bottom-up traversal of the syntax-tree of  $\psi_2$  to eliminate all the  $\wedge$ -symbols in  $\psi_2$ . The resulting formula  $\psi_3$  is obtained in time polynomial in the size of  $\psi_2$ . Furthermore,  $\psi_3$  is of the desired form  $\bigvee_{i \in I} \text{sph}_{\tau_i}(\bar{x})$  for an  $I \subseteq [\ell]$ . The overall time for constructing  $\psi_3$  and  $I$  is  $\text{poly}(\|\psi\|) \cdot 2^{(kd^{r+1})^{O(\|\sigma\|)}}$ . This completes the proof of Lemma 12.2.  $\square$

For evaluating a Boolean combination  $\psi(\bar{x})$  of sphere-formulas and Hanf-sentences on a given  $\sigma$ -db  $D$ , an obvious approach is to first consider every Hanf-sentence  $\chi$  that occurs in  $\psi$ , to check if  $D \models \chi$ , and to replace every occurrence of  $\chi$  in  $\psi$  with **true** (resp., **false**) if  $D \models \chi$  (resp.,  $D \not\models \chi$ ). The resulting formula  $\psi'(\bar{x})$  is then transformed into a disjunction  $\psi''(\bar{x}) := \bigvee_{i \in I} \text{sph}_{\tau_i}(\bar{x})$  by Lemma 12.2, and the query result  $\psi(D) = \psi''(D)$  is obtained as the union of the query results  $\text{sph}_{\tau_i}(D)$  for all  $i \in I$ .

While this works well in the static setting (i.e., without database updates), in the dynamic setting we have to take care of the fact that database updates might change the status of a Hanf-sentence  $\chi$  in  $\psi$ , i.e., an update operation might turn a database  $D$  with  $D \models \chi$  into a database  $D'$  with  $D' \not\models \chi$  (and vice versa). Consequently, the formula  $\psi''(\bar{x})$  that is equivalent to  $\psi(\bar{x})$  on  $D$  might be inequivalent to  $\psi(\bar{x})$  on  $D'$ .

To handle the dynamic setting correctly, at the end of each update step we will use the following lemma, which is an extension of Lemma 12.2 and is proved in a similar way.



### 12.1. Technical Lemmas on Types and Spheres Useful for Handling Non-Boolean Queries

**Lemma 12.3.** *Let  $\sigma$  be a schema. Let  $s \geq 0$  and let  $\chi_1, \dots, \chi_s$  be arbitrary formulas of schema  $\sigma$ . Let  $r \geq 0, k \geq 1, d \geq 2$  and let  $\mathcal{L}_r^{\sigma, d}(k) = \tau_1, \dots, \tau_\ell$ . Let  $\bar{x} = (x_1, \dots, x_k)$  be a list of  $k$  pairwise distinct variables. For every Boolean combination  $\psi(\bar{x})$  of the formulas  $\chi_1, \dots, \chi_s$  and of  $d$ -bounded sphere-formulas of radius at most  $r$  (over  $\sigma$ ), and for every  $J \subseteq [s]$  there is a set  $I \subseteq [\ell]$  such that*

$$\psi_J(\bar{x}) \equiv_d \bigvee_{i \in I} \text{sph}_{\tau_i}(\bar{x}),$$

where  $\psi_J$  is the formula obtained from  $\psi$  by replacing every occurrence of a formula  $\chi_j$  with **true** if  $j \in J$  and with **false** if  $j \notin J$  (for every  $j \in [s]$ ).

Given  $\psi$  and  $J$ , the set  $I$  can be computed in time  $\text{poly}(\|\psi\|) \cdot 2^{(kd^{r+1})^{O(\|\sigma\|)}}$ .

To evaluate a single sphere-formula  $\text{sph}_\tau(\bar{x})$  for a given  $r$ -type  $\tau$  with  $k$  centres (over  $\sigma$ ), it will be useful to decompose  $\tau$  into its connected components as follows. Let  $\tau = (T, \bar{t})$  with  $\bar{t} = (t_1, \dots, t_k)$ . Consider the Gaifman graph  $G^T$  of  $T$  and let  $C_1, \dots, C_c$  be the vertex sets of the  $c$  connected components of  $G^T$ . For each connected component  $C_j$  of  $G^T$ , let  $\bar{t}_j$  be the subsequence of  $\bar{t}$  consisting of all elements of  $\bar{t}$  that belong to  $C_j$ , and let  $k_j$  be the length of  $\bar{t}_j$ . Since  $(T, \bar{t})$  is an  $r$ -type with  $k$  centres, we have  $T = \mathcal{N}_r^T(\bar{t})$ , and thus  $c \leq k$  and  $k_j \geq 1$  for all  $j \in [c]$ . To avoid ambiguity, we make sure that the list  $C_1, \dots, C_c$  is sorted in such a way that for all  $j < j'$  we have  $i < i'$  for the smallest  $i$  with  $t_i \in C_j$  and the smallest  $i'$  with  $t_{i'} \in C_{j'}$ .

For each  $C_j$  consider the  $r$ -type with  $k_j$  centres  $\rho_j = (T[C_j], \bar{t}_j)$ . Let  $\nu_j$  be the unique integer such that  $\rho_j$  is isomorphic to the  $\nu_j$ -th element in the list  $\mathcal{L}_r^{\sigma, d}(k_j)$ , and let  $\tau_{j, \nu_j}$  be the  $\nu_j$ -th element in this list.

It is straightforward to see that the formula  $\text{sph}_\tau(\bar{x})$  is  $d$ -equivalent to the formula

$$\text{conn-sph}_\tau(\bar{x}) := \bigwedge_{j \in [c]} \text{sph}_{\tau_{j, \nu_j}}(\bar{x}_j) \wedge \bigwedge_{j \neq j'} \neg \text{dist}_{\leq 2r+1}^{k_j, k_{j'}}(\bar{x}_j, \bar{x}_{j'}), \quad (12.2)$$

where  $\bar{x}_j$  is the subsequence of  $\bar{x}$  obtained from  $\bar{x}$  in the same way as  $\bar{t}_j$  is obtained from  $\bar{t}$ , and  $\text{dist}_{\leq 2r+1}^{k_j, k_{j'}}(\bar{x}_j, \bar{x}_{j'})$  is a formula of schema  $\sigma$  which expresses that for some variable  $y$  in  $\bar{x}_j$  and some variable  $y'$  in  $\bar{x}_{j'}$  the distance between  $y$  and  $y'$  is  $\leq 2r+1$ . I.e., for  $\bar{a} = (a_1, \dots, a_{k_j})$  and  $\bar{b} = (b_1, \dots, b_{k_{j'}})$  we have  $(\bar{a}, \bar{b}) \in \text{dist}_{\leq 2r+1}^{k_j, k_{j'}}(D) \iff \text{dist}^D(\bar{a}; \bar{b}) \leq 2r+1$ , where

$$\text{dist}^D(\bar{a}; \bar{b}) \leq 2r+1 \text{ means that } \text{dist}^D(a_i, b_{i'}) \leq 2r+1 \text{ for some } i \in [k_j] \text{ and } i' \in [k_{j'}]. \quad (12.3)$$

Using the Lemmas 10.2 and 12.1, the following lemma is straightforward.

**Lemma 12.4.** *There is an algorithm which upon input of a schema  $\sigma$ , numbers  $r \geq 0, k \geq 1$  and  $d \geq 2$  and a  $r$ -type  $\tau$  with  $k$  centres (over  $\sigma$ ) computes the formula  $\text{conn-sph}_\tau(\bar{x})$ , along with the corresponding parameters  $c$  and  $k_j, \nu_j, \bar{x}_j, \tau_{j, \nu_j}$  for all  $j \in [c]$ .*

*The algorithm's runtime is  $2^{(kd^{r+1})^{O(\|\sigma\|)}}$ .*

## 12. Answering non-Boolean FO+MOD Queries Under Updates

We define the *signature* of  $\tau$  w.r.t.  $r$  to be the tuple  $\text{sgn}_r(\tau)$  built from the parameters  $c$  and  $(k_j, \nu_j, \{\mu \in [k] : x_\mu \text{ belongs to } \bar{x}_j\})_{j \in [c]}$  obtained from the above lemma. The signature  $\text{sgn}_r^D(\bar{a})$  of a tuple  $\bar{a}$  in a database  $D$  w.r.t. radius  $r$  is defined as  $\text{sgn}_r(\rho)$  for  $\rho := (\mathcal{N}_r^D(\bar{a}), \bar{a})$ . Note that  $\bar{a} \in \text{sph}_\tau(D) \iff \text{sgn}_r^D(\bar{a}) = \text{sgn}_r(\tau)$ .

### 12.2. Testing Non-Boolean FO+MOD Queries Under Updates

This section is devoted to the proof of the following theorem.

**Theorem 12.5.** *There is a dynamic algorithm that receives a schema  $\sigma$ , a degree bound  $d \geq 2$ , a  $k$ -ary FO+MOD[ $\sigma$ ]-query  $Q(\bar{x})$  (for some  $k \in \mathbb{N}$ ), and a  $\sigma$ -db  $D_0$  of degree  $\leq d$ , and computes within  $t_p = f(Q, d) \cdot \|D_0\|$  preprocessing time a data structure that can be updated in time  $t_u = f(Q, d)$  and allows to test for any input tuple  $\bar{a} \in \text{dom}^k$  whether  $\bar{a} \in \varphi(D)$  within testing time  $t_t = O(k^2)$ .*

For the proof, we use the lemmas provided in Section 12.1 and the following lemma.

**Lemma 12.6.** *There is a dynamic algorithm that receives a schema  $\sigma$ , a degree bound  $d \geq 2$ , numbers  $r \geq 0$  and  $k \geq 1$ , a  $r$ -type  $\tau$  with  $k$  centres (over  $\sigma$ ), and a  $\sigma$ -db  $D_0$  of degree  $\leq d$ , and computes within  $t_p = 2^{(kd^{r+1})^{O(\|\sigma\|)}} \cdot \|D_0\|$  preprocessing time a data structure that can be updated in time  $t_u = 2^{(kd^{r+1})^{O(\|\sigma\|)}}$  and allows to test for any input tuple  $\bar{a} \in \text{dom}^k$  whether  $\bar{a} \in \text{sph}_\tau(D)$  within testing time  $t_t = O(k^2)$ .*

*Proof.* The preprocessing routine starts by using Lemma 12.4 to compute the formula  $\text{conn-sph}_\tau(\bar{x})$ , along with the according parameters  $c$  and  $k_j, \nu_j, \bar{x}_j, \tau_{j, \nu_j}$  for each  $j \in [c]$ . This is done in time  $2^{(kd^{r+1})^{O(\|\sigma\|)}}$ . We let  $\text{sgn}_r(\tau)$  be the signature of  $\tau$  (defined directly after Lemma 12.4). Recall that  $\text{conn-sph}_\tau(\bar{x}) \equiv_d \text{sph}_\tau(\bar{x})$ , and recall from equation (12.2) the precise definition of the formula  $\text{conn-sph}_\tau(\bar{x})$ . Our data structure will store the following information on the database  $D$ :

- the set  $\Gamma$  of all tuples  $\bar{b} \in \text{adom}(D)^{k'}$  where  $k' \leq k$  and  $\mathcal{N}_r^D(\bar{b})$  is connected, and
- for every  $j \in [c]$  and every tuple  $\bar{b} \in \Gamma$  of arity  $k_j$ , the unique number  $\nu_{\bar{b}}$  such that  $\rho_{\bar{b}} := (\mathcal{N}_r^D(\bar{b}), \bar{b})$  is isomorphic to the  $\nu_{\bar{b}}$ -th element in the list  $\mathcal{L}_r^{\sigma, d}(k_j)$ .

We want to store this information in such a way that for any given tuple  $\bar{b} \in \text{dom}^{k'}$  it can be checked in time  $O(k)$  whether  $\bar{b} \in \Gamma$ . To ensure this, we use a  $k'$ -ary array  $\Gamma_{k'}$ <sup>1</sup> that is initialised to 0, and where during update operations the entry  $\Gamma_{k'}[\bar{b}]$  is set to 1 for all  $\bar{b} \in \Gamma$  of arity  $k'$ . In a similar way we can ensure that for any given  $j \in [c]$  and any  $\bar{b} \in \Gamma$  of arity  $k_j$ , the number  $\nu_{\bar{b}}$  can be looked up in time  $O(k)$ .

The **test** routine upon input of a tuple  $\bar{a} = (a_1, \dots, a_k)$  proceeds as follows.

First, we partition  $\bar{a}$  into  $\bar{a}_1, \dots, \bar{a}_{c'}$  (for  $c' \leq k$ ) such that  $C_j := \mathcal{N}_r^D(\bar{a}_j)$  for  $j \in [c']$  are the connected components of  $\mathcal{N}_r^D(\bar{a})$ . As in the definition of the formula

<sup>1</sup>This array requires nonlinear but polynomial space

$\text{conn-sph}_\tau(\bar{x})$ , we make sure that this list is sorted in such a way that for all  $j < j'$  we have  $i < i'$  for the smallest  $i$  with  $a_i \in C_j$  and the smallest  $i'$  with  $a_{i'} \in C_{j'}$ . All of this can be done in time  $O(k^2)$  by first constructing the graph  $H$  with vertex set  $[k]$  and where there is an edge between vertices  $i$  and  $j$  iff the tuple  $(a_i, a_j)$  belongs to  $\Gamma$  (i.e.,  $\mathcal{N}_r^D(a_i, a_j)$  is connected), and then computing the connected components of  $H$ .

Afterwards, for each  $j \in [c']$  we use time  $O(k)$  to look up the number  $\nu_{\bar{a}_j}$ . We then let  $\text{sgn}_r^D(\bar{a})$  be the tuple built from  $c'$  and  $(|\bar{a}_j|, \nu_{\bar{a}_j}, \{\mu \in [k] : a_\mu \text{ belongs to } \bar{a}_j\})_{j \in [c']}$ . It is straightforward to see that  $\bar{a} \in \text{conn-sph}_\tau(D)$  iff  $\text{sgn}_r^D(\bar{a}) = \text{sgn}_r(\tau)$ . Therefore, the **test** routine checks whether  $\text{sgn}_r^D(\bar{a}) = \text{sgn}_r(\tau)$  and outputs “yes” if this is the case and “no” otherwise. The entire time used by the **test** routine is  $t_t = O(k^2)$ .

To finish the proof of Lemma 12.6, we have to give further details on the **preprocess** routine and the **update** routine. The **preprocess** routine initialises  $\Gamma$  as the empty set  $\emptyset$  and then performs  $|D_0|$  update operations to insert all the tuples of  $D_0$  into the data structure. The **update** routine proceeds as follows.

Let  $D^-$  be the database before the update is received and let  $D^+$  be the database after the update has been performed. Let the update command be of the form **update**  $R(a_1, \dots, a_{\text{ar}(R)})$ . We let  $r' := r + (\text{ar}(R) - 1)(2r + 1)$ . All elements whose  $r'$ -neighbourhood might have changed belong to the set  $U := N_{r'}^{D'}(\bar{a})$ , where  $D' := D^+$  if the update command is **insert**  $R(\bar{a})$ , and  $D' := D^-$  if the update command is **delete**  $R(\bar{a})$ .

According to Lemma 10.2(d), all tuples  $\bar{b}$  that have to be inserted into or deleted from  $\Gamma$  are built from elements in  $U$ . To update the information stored in our data structure, we loop through all tuples of arity  $\leq k$  that are built from elements in  $U$ .

Using Lemma 10.2(a), we obtain that  $|U| \leq \text{ar}(R) \cdot d^{r'+1}$ . The number of candidate tuples  $\bar{b}$  built from elements in  $U$  is at most  $(\text{ar}(R) \cdot d^{r'+1})^{k+1}$ . Using the Lemmas 10.2 and 12.1, it is not difficult to see that the entire update time is at most  $t_u = 2^{(kd^{r+1})^{O(\|\sigma\|)}}$ . The initialisation time  $t_i$  is of the same form, and hence the preprocessing time is as claimed in the lemma. This completes the proof of Lemma 12.6.  $\square$

Using Lemma 12.6 and Lemma 12.3, we can show the following.

**Lemma 12.7.** *Let  $\sigma$  be a schema and let  $d \geq 2$  be a degree bound. Let  $s \geq 0$  and let  $\chi_1, \dots, \chi_s$  be arbitrary sentences of schema  $\sigma$ , and assume we have available for each  $j \in [s]$  a dynamic algorithm with initialisation time  $t'_i$  and update time  $t'_u$  that allows to check within answer time  $t'_{\text{ans}}$  whether or not  $D \models \chi_j$  for  $d$ -bounded  $\sigma$ -dbs  $D$ .*

*Then, there is a dynamic algorithm for  $d$ -bounded  $\sigma$ -dbs which receives as input numbers  $r \geq 0$  and  $k \geq 1$ , a tuple  $\bar{x} = (x_1, \dots, x_k)$  of pairwise distinct variables, and a Boolean combination  $\psi(\bar{x})$  of the sentences  $\chi_1, \dots, \chi_s$  and of  $d$ -bounded sphere-formulas of radius at most  $r$  (over  $\sigma$ ). Within initialisation time*

$$t_i = s(t'_i + t'_{\text{ans}}) + \text{poly}(\|\psi\|)2^{(kd^{r+1})^{O(\|\sigma\|)}} \quad (12.4)$$

*the algorithm builds a data structure that can be updated within time*

$$t_u = s(t'_u + t'_{\text{ans}}) + \text{poly}(\|\psi\|)2^{(kd^{r+1})^{O(\|\sigma\|)}} \quad (12.5)$$

## 12. Answering non-Boolean FO+MOD Queries Under Updates

and allows to test for any input tuple  $\bar{a} \in \text{adom}(D)^k$  whether  $\bar{a} \in \psi(D)$  within testing time

$$t_t = O(k^2). \quad (12.6)$$

*Proof.* We use Lemma 12.1 to compute the list  $\mathcal{L}_r^{\sigma,d}(k) = \tau_1, \dots, \tau_\ell$ . For each  $i \in [\ell]$ , we use the dynamic algorithm provided by Lemma 12.6 for  $\tau := \tau_i$ . Furthermore, for each  $j \in [s]$ , we use the dynamic algorithm provided by the lemma's assumption for checking whether or not  $D \models \chi_j$ . In addition to the components used by these dynamic algorithms, our data structure also stores

- the set  $J := \{j \in [s] : D \models \chi_j\}$ ,
- the particular set  $I \subseteq [\ell]$  provided by Lemma 12.3 for  $\psi(\bar{x})$  and  $J$ , and
- the set  $K = \{\text{sgn}_r(\tau_i) : i \in I\}$ , where for each type  $\tau$ ,  $\text{sgn}_r(\tau)$  is the signature of  $\tau$  defined directly after Lemma 12.4.

The **test** routine upon input of a tuple  $\bar{a} = (a_1, \dots, a_k)$  proceeds in the same way as in the proof of Lemma 12.6 to compute in time  $O(k^2)$  the signature  $\text{sgn}_r^D(\bar{a})$  of the tuple  $\bar{a}$ . For every  $i \in [\ell]$  we have  $\bar{a} \in \text{sph}_{\tau_i}(D) \iff \text{sgn}_r^D(\bar{a}) = \text{sgn}_r(\tau_i)$ . Thus,  $\bar{a} \in \varphi(D) \iff \text{sgn}_r^D(\bar{a}) \in K$ . Therefore, the **test** routine checks whether  $\text{sgn}_r^D(\bar{a}) \in K$  and outputs “yes” if this is the case and “no” otherwise. To ensure that this test can be done in time  $O(k^2)$ , we use an array construction for storing  $K$  (similar to the one for storing  $\Gamma$  in the proof of Lemma 12.6).

The **update** routine runs the update routines for all the used dynamic data structures. Afterwards, it recomputes  $J$  by calling the **answer** routine for  $\chi_j$  for all  $j \in [s]$ . Then, it uses Lemma 12.3 to recompute  $I$ . The set  $K$  is then recomputed by applying Lemma 12.4 for  $\tau := \tau_i$  for all  $i \in I$ . It is straightforward to verify that the initialisation time  $t_i$ , the update time  $t_u$ , and the testing time  $t_t$  are as claimed by the lemma.  $\square$

Theorem 12.5 is now obtained by combining Theorem 10.1, Lemma 12.7, and Theorem 11.1.

*Proof of Theorem 12.5.*

For  $k = 0$ , the theorem immediately follows from Theorem 11.1. Consider the case where  $k \geq 1$ . As in the proof of Theorem 11.1, we assume w.l.o.g. that all the symbols of  $\sigma$  occur in  $\varphi$ . We start the preprocessing routine by using Theorem 10.1 to transform  $\varphi(\bar{x})$  into a  $d$ -equivalent query  $\psi(\bar{x})$  in Hanf normal form; this takes time  $2^{d^{2^{O(\|\varphi\|)}}}$ . The formula  $\psi$  is a Boolean combination of  $d$ -bounded Hanf-sentences and sphere-formulas (over  $\sigma$ ) of locality radius at most  $r := 4^{\text{qr}(\varphi)}$ , and each sphere-formula is of arity at most  $k$ . Let  $\chi_1, \dots, \chi_s$  be the list of all Hanf-sentences that occur in  $\psi$ .

From Theorem 11.1 we have available for each  $j \in [s]$  a dynamic algorithm with initialisation time  $t'_i = O(\max_{j \in [s]} \|\chi_j\|)$  and update time  $t'_u = 2^{O(\|\sigma\|d^{2r+2})}$  that allows to check within answer time  $t'_{\text{ans}} = O(1)$  whether  $D \models \chi_j$  for  $d$ -bounded  $\sigma$ -dbs  $D$ . The proof of Theorem 12.5 therefore immediately follows from Lemma 12.7.  $\square$

## 12.3. Representing Databases by Coloured Graphs

To obtain dynamic algorithms for counting and enumerating query results, it will be convenient to work with a representation of databases by coloured graphs that is similar to the representation. The main advantage of this representation is that it unveils the combinatorial core of evaluating non-Boolean queries. In Theorem 12.8 we provide a general reduction from evaluating  $\text{FO}+\text{MOD}[\sigma]$ -queries to finding coloured independent sets in coloured graphs, which allows us to focus on this combinatorial graph problem when presenting the algorithms for counting (Section 12.4) and enumeration (Section 12.5).

For defining this representation, let us consider a fixed  $d$ -bounded  $r$ -type  $\tau$  with  $k$  centres (over a schema  $\sigma$ ). Use Lemma 12.4 to compute the formula  $\text{conn-sph}_\tau(\bar{x})$  (for  $\bar{x} = (x_1, \dots, x_k)$ ) and the according parameters  $c$  and  $k_j, \nu_j, \bar{x}_j, \tau_j, \nu_j$ , and let  $\text{sgn}_\tau(\tau)$  be the signature of  $\tau$ . To keep the notation simple, we assume w.l.o.g. that  $\bar{x}_1 = x_1, \dots, x_{k_1}, \bar{x}_2 = x_{k_1+1}, \dots, x_{k_1+k_2}$  etc.

Recall that  $\text{sph}_\tau(\bar{x})$  is  $d$ -equivalent to the formula

$$\text{conn-sph}_\tau(\bar{x}) \quad := \quad \bigwedge_{j \in [c]} \text{sph}_{\tau_j, \nu_j}(\bar{x}_j) \wedge \bigwedge_{j \neq j'} \neg \text{dist}_{\leq 2r+1}^{k_j, k_{j'}}(\bar{x}_j, \bar{x}_{j'}).$$

To count or enumerate the results of the formula  $\text{sph}_\tau(\bar{x})$  we represent the database  $D$  by a  $c$ -coloured graph  $\mathcal{G}_D$ . Here, a  $c$ -coloured graph  $\mathcal{G}$  is a database of the particular schema

$$\sigma_c \quad := \quad \{E, T_1, \dots, T_c\},$$

where  $E$  is a binary relation symbol and  $T_1, \dots, T_c$  are unary relation symbols. We define  $\mathcal{G}_D$  in such a way that the task of counting or enumerating the results of the query  $\text{sph}_\tau(\bar{x})$  on the database  $D$  can be reduced to counting or enumerating the results of the query

$$\varphi_c(z_1, \dots, z_c) \quad := \quad \bigwedge_{j \in [c]} T_j(z_j) \wedge \bigwedge_{j \neq j'} \neg E(z_j, z_{j'}) \quad (12.7)$$

on the  $c$ -coloured graph  $\mathcal{G}_D$ . The vertices of  $\mathcal{G}_D$  correspond to  $k_j$ -tuples over  $\text{adom}(D)$  (for some  $k_j \in \{k_1, \dots, k_c\} \subseteq [k]$ ) whose  $r$ -neighbourhood is connected; a vertex has colour  $T_j$  if its associated tuple  $\bar{a}$  is in  $\text{sph}_{\tau_j, \nu_j}(D)$ ; and an edge between two vertices indicates that  $\text{dist}^D(\bar{a}; \bar{b}) \leq 2r+1$ , for their associated tuples  $\bar{a}$  and  $\bar{b}$ . The following theorem allows to translate a dynamic algorithm for counting or enumerating the results of the query  $\varphi_c(z_1, \dots, z_c)$  on  $c$ -coloured graphs into a dynamic algorithm for counting or enumerating the result of an  $\text{FO}+\text{MOD}$ -query  $Q(\bar{x})$  on  $D$ .

**Theorem 12.8.** *Suppose for any  $d', c \in \mathbb{N}$  that the counting problem (the enumeration problem) for  $\varphi_c(\bar{z})$  on  $\sigma_c$ -dbs of degree at most  $d'$  can be solved by a dynamic algorithm with initialisation time  $t_i(c, d')$ , update time  $t_u(c, d')$  and counting time  $t_c(c, d')$  (delay  $t_d(c, d')$ ). Then for every schema  $\sigma$  and every  $d \geq 2$  the following holds. The counting problem (the enumeration problem) for  $k$ -ary  $\text{FO}+\text{MOD}[\sigma]$ -queries  $Q(\bar{x})$  on  $\sigma$ -dbs of degree at most  $d$  can be solved with counting time  $O(1)$  (delay  $O(\hat{t}_d + k)$ ), initialisation time  $\hat{t}_i \cdot f(Q, d)$  and update time*

## 12. Answering non-Boolean FO+MOD Queries Under Updates

- $(\hat{t}_u + \hat{t}_c) \cdot f(Q, d)$  for the counting problem and
- $\hat{t}_u \cdot f(Q, d)$  for the enumeration problem

where  $\hat{t}_x = \max_{c=1}^k t_x(c, d^{2^{O(\|\varphi\|)}})$  for  $t_x \in \{t_i, t_u, t_c, t_d\}$ .

The proof will be obtained as an easy consequence of Theorem 10.1, Theorem 11.1 and the following lemma.

**Lemma 12.9.** *Suppose that the counting problem (the enumeration problem) for  $\varphi_c(\bar{z})$  on  $\sigma_c$ -dbs of degree at most  $d'$  can be solved by a dynamic algorithm with initialisation time  $t_i(c, d')$ , update time  $t_u(c, d')$  and counting time  $t_c(c, d')$  (delay  $t_d(c, d')$ ).*

*Then for every schema  $\sigma$  and every  $d \geq 2$  the following holds. Let  $r \geq 0$ ,  $k \geq 1$ , and fix  $d' := d^{2k^2(2r+1)}$  and  $\tilde{t}_x := \max_{c=1}^k t_x(c, d')$  for  $t_x \in \{t_i, t_u, t_c, t_d\}$ .*

- (1) *Let  $\tau$  be a  $d$ -bounded  $r$ -type with  $k$  centres. The counting problem (the enumeration problem) for  $\text{sph}_\tau(\bar{x})$  on  $\sigma$ -dbs of degree at most  $d$  can be solved by a dynamic algorithm with counting time  $\tilde{t}_c$  (delay  $O(\tilde{t}_d + k)$ ), initialisation time  $\tilde{t}_i$ , and update time at most  $\tilde{t}_u d^{O(k^2 r + k\|\sigma\|)} + 2^{O(\|\sigma\|k^2 d^{2r+2})}$ .*
- (2) *Let  $s \geq 0$  and let  $\chi_1, \dots, \chi_s$  be arbitrary sentences of schema  $\sigma$ , and assume we have available for each  $j \in [s]$  a dynamic algorithm with initialisation time  $t'_i$  and update time  $t'_u$  that allows to check within answer time  $t'_{ans}$  whether or not  $D \models \chi_j$  for  $d$ -bounded  $\sigma$ -dbs  $D$ .*

*Let  $\bar{x} = (x_1, \dots, x_k)$  be a tuple of pairwise distinct variables, and let  $\psi(\bar{x})$  be a Boolean combination of the sentences  $\chi_1, \dots, \chi_s$  and of  $d$ -bounded sphere-formulas of radius at most  $r$  (over  $\sigma$ ).*

*Then, the counting problem (the enumeration problem) for  $\psi(\bar{x})$  on  $\sigma$ -dbs  $D$  of degree at most  $d$  can be solved by a dynamic algorithm with counting time  $O(1)$  (delay  $O(\tilde{t}_d + k)$ ), initialisation time  $s(t'_i + t'_{ans}) + 2^{(kd^{r+1})^{O(\|\sigma\|)}} (\text{poly}(\|\psi\|) + \tilde{t}_i)$ , and update time at most*

- $s(t'_u + t'_{ans}) + 2^{(k^2 d^{2r+2})^{O(\|\sigma\|)}} (\text{poly}(\|\psi\|) + \tilde{t}_c + \tilde{t}_u d^{O(k^2 r + k\|\sigma\|)})$  for the counting problem, and
- $s(t'_u + t'_{ans}) + 2^{(k^2 d^{2r+2})^{O(\|\sigma\|)}} (\text{poly}(\|\psi\|) + \tilde{t}_u d^{O(k^2 r + k\|\sigma\|)})$  for the enumeration problem.

*Proof.* We prove part (1) by a reduction from  $\text{conn-sph}_\tau(\bar{x})$  to  $\varphi_c$ . We use the notation introduced at the beginning of Section 12.3, and we let  $\tau_j := \tau_{j, \nu_j}$  for every  $j \in [c]$ . For a  $\sigma$ -db  $D$  of degree at most  $d$  we let  $\mathcal{G}_D$  be the  $\sigma_c$ -db with

$$\begin{aligned} T_j^{\mathcal{G}_D} &:= \{ v_{\bar{a}} : \bar{a} \in \text{adom}(D)^{k_j} \text{ with } (\mathcal{N}_r^D(\bar{a}), \bar{a}) \cong \tau_j \}, \quad \text{for all } j \in [c], \text{ and} \\ E^{\mathcal{G}_D} &:= \left\{ (v_{\bar{a}}, v_{\bar{b}}) \in V^2 : \text{dist}^D(\bar{a}; \bar{b}) \leq 2r+1 \right\}, \end{aligned}$$

where  $V := \bigcup_{j \in [c]} T_j^{\mathcal{G}_D}$ . We will shortly write  $E$  and  $T_j$  instead of  $E^{\mathcal{G}_D}$  and  $T_j^{\mathcal{G}_D}$ .

### 12.3. Representing Databases by Coloured Graphs

Using Lemma 10.2 and the fact that  $\tau_j$  is connected we obtain that  $(v_{\bar{a}}, v_{\bar{b}}) \in E$  iff  $\mathcal{N}_r^D(\bar{a}, \bar{b})$  is connected. If  $\mathcal{N}_r^D(\bar{a}, \bar{b})$  is connected, then by Lemma 10.2(d)  $\bar{b} \in (N_{r+(|\bar{a}|+|\bar{b}|-1)(2r+1)}^D(a_1))^{\bar{b}}$ . It follows that the degree of  $\mathcal{G}_D$  is bounded by  $d^{2k^2(2r+1)}$ . Furthermore, by the definition of  $\mathcal{G}_D$  and  $\varphi_c$  we get that  $(\bar{a}_1, \dots, \bar{a}_c) \in \mathbf{sph}_\tau(D) \iff (v_{\bar{a}_1}, \dots, v_{\bar{a}_c}) \in \varphi_c(\mathcal{G}_D)$ , for all tuples  $\bar{a}_1, \dots, \bar{a}_c$  where  $\bar{a}_j$  has arity  $k_j$  for each  $j \in [c]$ . As a consequence,  $|\mathbf{sph}_\tau(D)| = |\varphi_c(\mathcal{G}_D)|$ , and we can therefore use the **count** routine for  $\varphi_c$  on  $\mathcal{G}_D$  to count the number of tuples in  $\mathbf{sph}_\tau(D)$ . Furthermore, by annotating every vertex  $v_{\bar{a}}$  with its tuple  $\bar{a}$ , we can translate every tuple  $(v_{\bar{a}_1}, \dots, v_{\bar{a}_c}) \in \varphi_c(\mathcal{G}_D)$  to  $(\bar{a}_1, \dots, \bar{a}_c)$  in time  $O(k)$ . Therefore, given an **enumerate** routine for  $\varphi_c(\mathcal{G}_D)$  with delay  $t_d$  we can produce an enumeration of  $\mathbf{sph}_\tau(D)$  with delay  $O(t_d + k)$ .

It remains to show how to construct and maintain  $\mathcal{G}_D$  when the database  $D$  is updated. As initialisation for the empty database  $D_\emptyset$  we just perform the **init** routine of the dynamic algorithm for  $\varphi_c(\bar{z})$  on  $\sigma_c$ -dbs of degree at most  $d'$ . The **update** routine of the dynamic algorithm for  $\mathbf{sph}_\tau(\bar{x})$  on  $\sigma$ -dbs of degree at most  $d$  is provided by the following claim.

**Claim 12.10.** *If  $D^+$  is obtained from  $D^-$  by one update step, then  $\mathcal{G}_{D^+}$  can be obtained from  $\mathcal{G}_{D^-}$  by  $d^{O(k^2 r + k \|\sigma\|)}$  update steps and additional computing time  $2^{O(\|\sigma\| k^2 d^{2r+2})}$ .*

*Proof.* Let the update command be of the form **update**  $R(a_1, \dots, a_{\text{ar}(R)})$  with  $\bar{a} = (a_1, \dots, a_{\text{ar}(R)})$ . Let  $D' \in \{D^-, D^+\}$  be the database whose relation  $R$  contains the tuple  $\bar{a}$  (either before deletion or after insertion). Let  $r' := r + (k-1)(2r+1)$  and note that all elements in the active domain whose  $r'$ -neighbourhood in the database might have changed, belong to the set  $U := N_{r'}^{D'}(\bar{a})$ .

For every  $j \in [c]$  and every tuple  $\bar{b}$  of arity at most  $k$  of elements in  $U$ , we check whether the  $r$ -type  $(\mathcal{N}_r^{D^+}(\bar{b}), \bar{b})$  of  $\bar{b}$  is isomorphic to  $\tau_j$ . Depending on the outcome of this test, we include or exclude  $v_{\bar{b}}$  from the relation  $T_j$ . Note that it indeed suffices to consider the tuples  $\bar{b}$  built from elements in  $U$ : The  $r$ -type of some tuple  $\bar{b}$  is changed by the update command only if  $N_{r'}^{D'}(\bar{b})$  contains some element from  $\bar{a}$ . Furthermore, we only have to consider tuples  $\bar{b}$  whose  $r$ -neighbourhood  $\mathcal{N}_r^{D'}(\bar{b})$  is connected. Using Lemma 10.2(d), we therefore obtain that each component of  $\bar{b}$  belongs to  $N_{r'}^{D'}(\bar{a}) = U$ .

Afterwards, we update the coloured graph's edge relation  $E$ . There is an edge  $(v_{\bar{b}}, v_{\bar{b}'}) \in E^{D^+}$ , if and only if there are  $j, j' \in [c]$  such that  $v_{\bar{b}} \in C_j$ ,  $v_{\bar{b}'} \in C_{j'}$  and  $\text{dist}^{D^+}(\bar{b}; \bar{b}') \leq 2r+1$ . Note that if  $\text{dist}^{D^+}(\bar{b}; \bar{b}') \leq 2r+1$ , then there is some component  $b_i$  in  $\bar{b}$  and some component  $b_{i'}$  in  $\bar{b}'$  such that  $\text{dist}^{D^+}(b_i, b_{i'}) \leq 2r+1$ ; and in case that  $v_{\bar{b}'} \in C_{j'}$ , we know that  $\mathcal{N}_r^{D^+}(\bar{b}')$  is connected, and hence every component of the tuple  $\bar{b}'$  belongs to  $N_{r+(k-1)(2r+1)}^{D^+}(b_{i'}) \subseteq N_{r+k(2r+1)}^{D^+}(b_i) \subseteq N_{r+k(2r+1)}^{D^+}(\bar{b})$ . Moreover, for correctly updating the edge relation  $E$  it suffices to consider only those pairs of tuples  $\bar{b}$  and  $\bar{b}'$  where for at least one of the two tuples, at least one component belongs to  $N_{r'}^{D'}(\bar{a})$ , and hence all components belong to  $N_{r'}^{D'}(\bar{a})$ , since these are the tuples where the  $r$ -neighbourhood or the condition  $\text{dist}^D(\bar{b}; \bar{b}') \leq 2r+1$  might be affected by the database update. Overall, the algorithm proceeds as follows: We compute for all tuples  $\bar{b}$  of arity at most  $k$  in  $N_{r'}^{D'}(\bar{a})$ , all tuples  $\bar{b}'$  of arity at most  $k$  in  $N_{r+k(2r+1)}^{D'}(\bar{b})$

## 12. Answering non-Boolean FO+MOD Queries Under Updates

(later on, we will call these tuples *candidate tuples*) and check whether (1) there is a  $j \in [c]$  such that  $v_{\bar{b}} \in T_j$ , (2) there is a  $j' \in [c]$  such that  $v_{\bar{b}'} \in T_{j'}$ , and (3)  $\text{dist}^{D^+}(\bar{b}; \bar{b}') \leq 2r+1$ . If all three checks return the result “yes”, then we insert the tuple  $(v_{\bar{b}}, v_{\bar{b}'})$  into  $E$ , otherwise we remove it from  $E$ .

It remains to analyse the runtime of the described update procedure. By Lemma 10.2,  $|U| \leq \text{ar}(R)d^{r'+1} \leq \|\sigma\|d^{r+(k-1)(2r+1)+1} \leq d^{O(kr+\lg \|\sigma\|)} \leq d^{O(kr+\|\sigma\|)}$ . Furthermore,  $U$  can be computed in time  $(\text{ar}(R)d^{r'+1})^{O(\|\sigma\|)} \leq d^{O(kr\|\sigma\|+\|\sigma\|^2)}$ . The number of tuples  $\bar{b}$  that we have to consider is at most  $\sum_{i=1}^k |U|^i \leq |U|^{k+1} \leq d^{O(k^2r+k\|\sigma\|)}$ .

For each such  $\bar{b}$  we use Lemma 10.2(e) to check in time  $2^{O(\|\sigma\|k^2d^{2r+2})}$  whether the  $r$ -type of  $\bar{b}$  in  $D^+$  is isomorphic to  $\tau_j$ , for some  $j \in [c]$ . In summary, for updating the sets  $T_1, \dots, T_c$  we use at most  $c|U|^{k+1} \leq d^{O(k^2r+k\|\sigma\|)}$  calls of the **update** routine of the dynamic algorithm on coloured graphs, and in addition to that we use computation time at most  $2^{O(\|\sigma\|k^2d^{2r+2})}$ .

To update the edge relation, we compute for each of the  $d^{O(k^2r+k\|\sigma\|)}$  tuples  $\bar{b}$  a list of all its candidate tuples  $\bar{b}'$ ; using Lemma 10.2 this can be done in time  $(\text{ar}(R)d^{r+k(2r+1)+1})^{O(\|\sigma\|)} \leq d^{O(kr\|\sigma\|+\|\sigma\|^2)}$ . By Lemma 10.2, the number of candidate tuples is  $\leq d^{O(k^2r+k\|\sigma\|)}$ . By Lemma 10.2(e), it takes time  $2^{O(\|\sigma\|k^2d^{2r+2})}$  to check if  $(\mathcal{N}_r^{D^+}(\bar{b}), \bar{b})$  is isomorphic to  $\tau_j$  and  $(\mathcal{N}_r^{D^+}(\bar{b}'), \bar{b}')$  is isomorphic to  $\tau_{j'}$ , for some  $j, j' \in [c]$ . We can use and maintain an additional array that allows us to check, for any  $a_i$  and  $b_j$ , in constant time whether  $\text{dist}^D(a_i, b_j) \leq 2r+1$ . Overall, we obtain that also the edge relation  $E$  can be updated by at most  $d^{O(k^2r+k\|\sigma\|)}$  calls of the **update** routine of the dynamic algorithm on coloured graphs and additional computation time at most  $2^{O(\|\sigma\|k^2d^{2r+2})}$ .

This completes the proof of Claim 12.10.  $\square$

Finally, the **preprocess** routine of the dynamic algorithm for  $\text{sph}_\tau(\bar{x})$  proceeds in the obvious way by first calling the **init** routine for  $D_\emptyset$  and then performing  $|D_\emptyset|$  update steps to insert all the tuples of  $D_\emptyset$  into the data structure. This completes the proof of part (1) of Lemma 12.9.

We now turn to the proof of part (2) of Lemma 12.9. We use Lemma 12.1 to compute the list  $\mathcal{L}_r^{\sigma,d}(k) = \tau_1, \dots, \tau_\ell$  within time  $2^{(kd^{r+1})^{O(\|\sigma\|)}}$ . For each  $i \in [\ell]$ , we use the dynamic algorithm for  $\text{sph}_{\tau_i}(\bar{x})$  provided from the lemma's part (1). Furthermore, for each  $j \in [s]$ , we use the dynamic algorithm for answering whether or not  $D \models \chi_j$ , provided by the assumption of part (2) of Lemma 12.9. In addition to the components used by these dynamic algorithms, our data structure also stores

- the set  $J := \{j \in [s] : D \models \chi_j\}$  and
- the particular set  $I \subseteq [\ell]$  provided by Lemma 12.3 for  $\psi(\bar{x})$  and  $J$ .

For the case that we want to solve the counting problem, our data structure also stores

- the cardinality  $n = |\varphi(D)|$  of the query result.



### 12.3. Representing Databases by Coloured Graphs

The **count** routine simply outputs the value  $n$  in time  $O(1)$ .

The **enumerate** routine runs the **enumerate** routine on  $\text{sph}_{\tau_i}(D)$  for every  $i \in I$ . Note that this enumerates, without repetition, all tuples in  $\varphi(D)$ , because by Lemma 12.3,  $\varphi(D)$  is the union of the sets  $\text{sph}_{\tau_i}(D)$  for all  $i \in I$ , and this is a union of pairwise disjoint sets.

The **update** routine runs the update routines for all used dynamic data structures. Afterwards, it recomputes  $J$  by calling the **answer** routine for  $\chi_j$  for all  $j \in [s]$ . Then, it uses Lemma 12.3 to recompute  $I$ . For the case that we want to solve the counting problem, we afterwards recompute the number  $n$  by letting  $n = \sum_{i \in I} n_i$ , where  $n_i$  is the result of the **count** routine for  $\tau_i$ .

By using the statement of part (1) and the assumptions of part (2) of the lemma, it is straightforward to verify that the initialisation time, the update time, and the counting time (the delay) are as claimed by the lemma.  $\square$

Theorem 12.8 is now obtained by combining Theorem 10.1, part (2) of Lemma 12.9, and Theorem 11.1.

*Proof of Theorem 12.8.*

For  $k = 0$ , the result follows immediately from Theorem 11.1. Consider the case where  $k \geq 1$ . W.l.o.g. we assume that all the symbols of  $\sigma$  occur in  $Q$  (otherwise, we remove from  $\sigma$  all symbols that do not occur in  $Q$ ). We start the preprocessing routine by using Theorem 10.1 to transform  $\varphi(\bar{x})$  into a  $d$ -equivalent query  $\psi(\bar{x})$  in Hanf normal form; this takes time  $2^{d^{2^{O(\|Q\|)}}}$ . The formula  $\psi$  is a Boolean combination of  $d$ -bounded Hanf-sentences and sphere-formulas (over  $\sigma$ ) of locality radius at most  $r := 4^{\text{qr}(\varphi)}$ , and each sphere-formula is of arity at most  $k$ . Note that for  $d' := d^{2k^2(2r+1)}$  as used in the lemma's part (1), it holds that  $d' = d^{2^{O(\|\varphi\|)}}$ . Let  $\chi_1, \dots, \chi_s$  be the list of all Hanf-sentences that occur in  $\psi$ , and note that  $s \leq 2^{d^{2^{O(\|\varphi\|)}}}$ .

From Theorem 11.1 we have available for each  $j \in [s]$  a dynamic algorithm with initialisation time  $t'_i = O(\max_{j \in [s]} \|\chi_j\|)$  and update time  $t'_u = 2^{O(\|\sigma\|d^{2r+2})}$  that allows to check within answer time  $t'_{\text{ans}} = O(1)$  whether  $D \models \chi_j$  for  $d$ -bounded  $\sigma$ -dbs  $D$ .

Applying part (2) of Lemma 12.9, we obtain a dynamic algorithm that solves the counting problem (the enumeration problem) for  $\varphi(\bar{x})$  on  $\sigma$ -dbs of degree at most  $d$  with counting time  $O(1)$  (delay  $O(\tilde{t}_d + k) = O(\hat{t}_d + k)$ ), initialisation time

$$s(t'_i + t'_{\text{ans}}) + 2^{(kd^{r+1})^{O(\|\sigma\|)}} (\text{poly}(\|\psi\|) + \tilde{t}_i) \leq \hat{t}_i \cdot 2^{d^{2^{O(\|\varphi\|)}}}$$

and update time at most

$$\begin{aligned} & s(t'_u + t'_{\text{ans}}) + 2^{(k^2 d^{2r+2})^{O(\|\sigma\|)}} (\text{poly}(\|\psi\|) + \tilde{t}_c + \tilde{t}_u d^{O(k^2 r + k\|\sigma\|)}) \\ & \leq (\hat{t}_u + \hat{t}_c) \cdot 2^{d^{2^{O(\|\varphi\|)}}} \end{aligned}$$

when dealing with the counting problem and update time at most

$$s(t'_u + t'_{\text{ans}}) + 2^{(k^2 d^{2r+2})^{O(\|\sigma\|)}} (\text{poly}(\|\psi\|) + \tilde{t}_u d^{O(k^2 r + k\|\sigma\|)}) \leq \hat{t}_u \cdot 2^{d^{2^{O(\|\varphi\|)}}}$$

when dealing with the enumeration problem.  $\square$

## 12.4. Counting Results of FO+MOD Queries Under Updates

This section is devoted to the proof of the following theorem.

**Theorem 12.11.** *There is a dynamic algorithm that receives a schema  $\sigma$ , a degree bound  $d \geq 2$ , a  $k$ -ary FO+MOD[ $\sigma$ ]-query  $Q(\bar{x})$  (for some  $k \in \mathbb{N}$ ) and a  $\sigma$ -db  $D_0$  of degree  $\leq d$  and computes within  $t_p = f(Q, d) \cdot \|D_0\|$  preprocessing time a data structure that can be updated in time  $t_u = f(Q, d)$  and allows to return the cardinality  $|Q(D)|$  of the query result within time  $O(1)$ .*

The theorem follows immediately from Theorem 12.8 and the following dynamic counting algorithm for the query  $\varphi_c(\bar{z})$ .

**Lemma 12.12.** *There is a dynamic algorithm that receives a number  $c \geq 1$ , a degree bound  $d \geq 2$  and a  $\sigma_c$ -db  $\mathcal{G}_0$  of degree  $\leq d$ , and computes  $|\varphi_c(\mathcal{G})|$  with  $d^{O(c^2)}$  initialisation time,  $O(1)$  counting time, and  $d^{O(c^2)}$  update time.*

*Proof.* Recall from (12.7) that  $\varphi_c(z_1, \dots, z_c) = \bigwedge_{i \in [c]} T_i(z_i) \wedge \bigwedge_{j \neq j'} \neg E(z_j, z_{j'})$ . For all  $j, j' \in [c]$  with  $j \neq j'$  consider the formula  $\theta_{j,j'}(z_1, \dots, z_c) := E(z_j, z_{j'}) \wedge \bigwedge_{i \in [c]} T_i(z_i)$ . Furthermore, let  $\alpha(z_1, \dots, z_c) := \bigwedge_{i \in [c]} T_i(z_i)$ . Clearly, for every  $\sigma_c$ -db  $\mathcal{G}$  we have

$$\begin{aligned} \alpha(\mathcal{G}) &= T_1^{\mathcal{G}} \times \dots \times T_c^{\mathcal{G}}, \\ \varphi_c(\mathcal{G}) &= \alpha(\mathcal{G}) \setminus \left( \bigcup_{j \neq j'} \theta_{j,j'}(\mathcal{G}) \right), \\ \text{and hence, } |\varphi_c(\mathcal{G})| &= |\alpha(\mathcal{G})| - \left| \bigcup_{j \neq j'} \theta_{j,j'}(\mathcal{G}) \right|. \end{aligned}$$

By the *inclusion-exclusion principle* (see Equation 2.2) we obtain for

$$J := \{(j, j') : j, j' \in [c], j \neq j'\}$$

that

$$\left| \bigcup_{j \neq j'} \theta_{j,j'}(\mathcal{G}) \right| = \sum_{\emptyset \neq K \subseteq J} (-1)^{|K|-1} \left| \bigcap_{(j,j') \in K} \theta_{j,j'}(\mathcal{G}) \right| = \sum_{\emptyset \neq K \subseteq J} (-1)^{|K|-1} |\varphi_K(\mathcal{G})|$$

for the formula  $\varphi_K(z_1, \dots, z_c) := \bigwedge_{i \in [c]} T_i(z_i) \wedge \bigwedge_{(j,j') \in K} E(z_j, z_{j'})$ . Our data structure stores the following values:

- $|T_i^{\mathcal{G}}|$ , for each  $i \in [c]$ , and  $n_1 := |\alpha(\mathcal{G})| = \prod_{i \in [c]} |T_i^{\mathcal{G}}|$ ,
- $|\varphi_K(\mathcal{G})|$ , for each  $K \subseteq J$  with  $K \neq \emptyset$ , and
- $n_2 := \sum_{\emptyset \neq K \subseteq J} (-1)^{|K|-1} |\varphi_K(\mathcal{G})|$  and  $n_3 := n_1 - n_2$ .

#### 12.4. Counting Results of FO+MOD Queries Under Updates

Note that  $n_3 = |\varphi_c(\mathcal{G})|$  is the desired size of the query result. Therefore, the **count** routine can answer in time  $O(1)$  by just outputting the number  $n_3$ .

It remains to show how these values can be initialised and updated during updates of  $\mathcal{G}$ . The initialisation for the empty graph initialises all the values to 0. In the **update** routine, the values for  $|T_i^{\mathcal{G}}|$  and  $n_1$  can be updated in a straightforward way (using time  $O(c)$ ). For each  $K \subseteq J$ , the update of  $|\varphi_K(\mathcal{G})|$  is provided within time  $d^{O(c^2)}$  by the following Claim 12.13.

**Claim 12.13.** *For every  $K \subseteq J$ , the cardinality  $|\varphi_K(\mathcal{G})|$  of a  $\sigma_c$ -db  $\mathcal{G}$  of degree at most  $d$  can be updated within time  $d^{O(c^2)}$  after  $d^{O(c^2)} \cdot |\mathcal{G}_0|$  preprocessing time.*

*Proof.* Consider the directed graph  $H := (V, K)$  with vertex set  $V := [c]$  and edge set  $K$ . Decompose the Gaifman graph of  $H$  into its connected components. Let  $V_1, \dots, V_s$  be the connected components (for a suitable  $s \leq c$ ). For each  $i \in [s]$  let  $H_i := H[V_i]$  be the induced subgraph of  $H$  on  $V_i$ . We write  $K_i$  to denote the set of edges of  $H_i$ . For every  $i \in [s]$  let  $\ell_i = |V_i|$  and let  $t(i, 1) < t(i, 2) < \dots < t(i, \ell_i)$  be the ordered list of the vertices in  $V_i$ . Consider the query

$$\varphi_{K_i}(z_{t(i,1)}, \dots, z_{t(i,\ell_i)}) := \bigwedge_{j \in V_i} T_j(z_j) \wedge \bigwedge_{(j,j') \in K_i} E(z_j, z_{j'}). \quad (12.8)$$

Note that  $\varphi_K$  is the conjunction of the formulas  $\varphi_{K_i}$  for all  $i \in [s]$ . Since the variables of the formulas  $\varphi_{K_i}$  for  $i \in [s]$  are pairwise disjoint, we have  $|\varphi_K(\mathcal{G})| = |\varphi_{K_1}(\mathcal{G})| \times \dots \times |\varphi_{K_s}(\mathcal{G})|$  (modulo permutations of the tuples), and thus  $|\varphi_K(\mathcal{G})| = \prod_{i \in [s]} |\varphi_{K_i}(\mathcal{G})|$ .

For each  $i \in [s]$ , the value  $|\varphi_{K_i}(\mathcal{G})|$  can be computed as follows. For every  $v \in \text{adom}(\mathcal{G})$  we consider the set  $S_i^v := \{(w_{t(i,1)}, \dots, w_{t(i,\ell_i)}) \in \varphi_{K_i}(\mathcal{G}) : w_{t(i,1)} = v\}$ . Since the Gaifman graph of  $H_i$  is connected and has  $\ell_i$  nodes, it follows that each component of every tuple in  $S_i^v$  is contained in the  $(\ell_i - 1)$ -neighbourhood of  $v$  in  $\mathcal{G}$ , and this neighbourhood contains at most  $d^{\ell_i}$  elements. Therefore,  $|S_i^v| \leq d^{(\ell_i)^2}$ , and using breadth-first search starting from  $v$ , the set  $S_i^v$  can be computed in time  $d^{O(c^2)}$ . Note that  $\varphi_{K_i}(\mathcal{G})$  is the disjoint union of the sets  $S_i^v$  for all  $v \in \text{adom}(\mathcal{G})$ . Therefore,  $|\varphi_{K_i}(\mathcal{G})| = \sum_{v \in \text{adom}(\mathcal{G})} |S_i^v|$ .

In our data structure we store for every  $i \in [s]$  and every  $v \in \text{adom}(\mathcal{G})$  the number  $\mu_{i,v} = |S_i^v|$ . Moreover, for every  $i \in [s]$  we store the sum  $\mu_i = \sum_{v \in \text{adom}(\mathcal{G})} \mu_{i,v} = |\varphi_{K_i}(\mathcal{G})|$ .

The initialisation for the empty  $\sigma_c$ -db  $\mathcal{G}_0$  sets all these values to 0. Whenever the colour of a vertex of  $\mathcal{G}$  is updated or an edge is inserted or deleted, we update all affected numbers accordingly. Note that a number  $\mu_{i,v}$  changes only if  $v$  is in the  $(c - 1)$ -neighbourhood around the updated edge or vertex in the graph  $\mathcal{G}$ . Hence, for at most  $2d^c$  vertices  $v$ , the numbers  $\mu_{i,v}$  are affected by an update, and each of them can be updated in time  $d^{O(c^2)}$ . Moreover, for each  $i \in [s]$ , the sum  $\mu_i$  can be updated in time  $O(d^c)$  by subtracting the old value of  $\mu_{i,v}$  and adding the new value of  $\mu_{i,v}$  for each of the at most  $2d^c$  relevant vertices  $v$ . Finally, it takes time  $O(c)$  to compute the updated value  $|\varphi_K(\mathcal{G})| = \prod_{i \in [s]} \mu_i$ . The overall time used to produce the update is  $d^{O(c^2)}$ .  $\square$

## 12. Answering non-Boolean FO+MOD Queries Under Updates

Once we have available the updated numbers  $|\varphi_K(\mathcal{G})|$  for all  $K \subseteq J$ , the value  $n_2$  can be computed in time  $O(|2^J|) \leq 2^{O(c^2)}$ . And  $n_3$  is then obtained in time  $O(1)$ . Altogether, performing the **update** routine takes time at most  $d^{O(c^2)}$ . The **preprocess** routine initialises all values for the empty graph and then uses  $|\mathcal{G}_0|$  update steps to insert all the tuples of  $\mathcal{G}_0$  into the data structure. This completes the proof of Lemma 12.12.  $\square$

### 12.5. Enumerating Results of FO+MOD Queries Under Updates

In this section we prove — and afterwards improve — the following theorem.

**Theorem 12.14.** *There is a dynamic algorithm that receives a schema  $\sigma$ , a degree bound  $d \geq 2$ , a  $k$ -ary FO+MOD[ $\sigma$ ]-query  $Q(\bar{x})$  (for some  $k \in \mathbb{N}$ ) and a  $\sigma$ -db  $D_0$  of degree  $\leq d$ , and computes within  $t_p = f(Q, d) \cdot \|D_0\|$  preprocessing time a data structure that can be updated in time  $t_u = f(Q, d)$  and allows to enumerate  $Q(D)$  with  $d^{2^{O(\|\varphi\|)}}$  delay.*

The theorem follows immediately from Theorem 12.8 and the following dynamic enumeration algorithm for the query  $\varphi_c(\bar{z})$ .

**Lemma 12.15.** *There is a dynamic algorithm that receives a number  $c \geq 1$ , a degree bound  $d \geq 2$  and a  $\sigma_c$ -db  $\mathcal{G}_0$  of degree  $\leq d$ , and computes within  $t_p = d^{\text{poly}(c)} \cdot |\mathcal{G}_0|$  preprocessing time a data structure that can be updated in time  $d^{\text{poly}(c)}$  and allows to enumerate the query result  $\varphi_c(\mathcal{G})$  with  $O(c^3 d)$  delay.*

*Proof.* For a  $\sigma_c$ -db  $\mathcal{G}$  and a vertex  $v \in \text{adom}(\mathcal{G})$  we let  $N^\mathcal{G}(v)$  be the set of all neighbours of  $v$  in  $\mathcal{G}$ . I.e.,  $N^\mathcal{G}(v)$  is the set of all  $w \in \text{adom}(\mathcal{G})$  such that  $(v, w)$  or  $(w, v)$  belongs to  $E^\mathcal{G}$ .

The underlying idea of the enumeration procedure is the following greedy strategy. We cycle through all vertices  $u_1 \in T_1^\mathcal{G}$ ,  $u_2 \in T_2^\mathcal{G} \setminus N^\mathcal{G}(u_1)$ ,  $u_3 \in T_3^\mathcal{G} \setminus (N^\mathcal{G}(u_1) \cup N^\mathcal{G}(u_2))$ ,  $\dots$ ,  $u_c \in T_c^\mathcal{G} \setminus \bigcup_{i \leq c-1} N^\mathcal{G}(u_i)$  and output  $(u_1, \dots, u_c)$ . This strategy does not yet lead to a constant delay enumeration, as there might be vertex tuples  $(u_1, \dots, u_i)$  (for  $i < c$ ) that do extend to an output tuple  $(u_1, \dots, u_c)$ , but where many possible extensions are checked before this output tuple is encountered. We now show how to overcome this problem and describe an enumeration procedure with  $O(c^3 d)$  delay and update time  $d^{\text{poly}(c)}$ .

Note that for every  $J \subseteq [c]$  we have  $|\bigcup_{j \in J} N^\mathcal{G}(u_j)| \leq cd$ . Hence, if a set  $T_i^\mathcal{G}$  contains more than  $cd$  elements, we know that *every* considered tuple has an extension  $u_i \in T_i^\mathcal{G}$  that is not a neighbour of any vertex in the tuple. Let  $I := \{i \in [c] : |T_i^\mathcal{G}| \leq cd\}$  be the set of *small* colour classes in  $\mathcal{G}$  and to simplify the presentation we assume without loss of generality that  $I = \{1, \dots, s\}$ . In our data structure we store the current index set  $I$  and the set

$$\mathcal{S} := \{ (u_1, \dots, u_s) \in T_1^\mathcal{G} \times \dots \times T_s^\mathcal{G} : (u_j, u_{j'}) \notin E^\mathcal{G}, \text{ for all } j \neq j' \} \quad (12.9)$$

### 12.5. Enumerating Results of FO+MOD Queries Under Updates

of tuples on the small colours. Note that a tuple  $(u_1, \dots, u_s) \in T_1^{\mathcal{G}} \times \dots \times T_s^{\mathcal{G}}$  extends to an output tuple  $(u_1, \dots, u_c) \in \varphi_c(\mathcal{G})$  if and only if it is contained in  $\mathcal{S}$ . We store the current sizes of all colours and this enables us to keep the set  $I$  of small colours updated. Moreover, as  $|\mathcal{S}| \leq (cd)^c$ , we can update the set  $\mathcal{S}$  in time  $d^{\text{poly}(c)}$  after every update by a brute-force approach. The enumeration procedure is given in Algorithm 18.

---

**Algorithm 18** Enumeration procedure with delay  $O(c^3d)$ .

---

```

1: for all  $(u_1, \dots, u_s) \in \mathcal{S}$  do
2:   ENUM( $u_1, \dots, u_s$ ).
3: Output the end-of-enumeration message EOE.
4:
5: function ENUM( $u_1, \dots, u_i$ )
6:   if  $i = c$  then
7:     Output the tuple  $(u_1, \dots, u_c)$ .
8:   else
9:     for all  $u_{i+1} \in T_{i+1}^{\mathcal{G}}$  do
10:      if  $u_{i+1} \notin \bigcup_{j=1}^i N^{\mathcal{G}}(u_j)$  then ENUM( $u_1, \dots, u_i, u_{i+1}$ ).

```

---

It is straightforward to see that this procedure enumerates  $\varphi_c(\mathcal{G})$ . Let us analyse the delay. Since for all  $i > s$  we have  $|T_i^{\mathcal{G}}| > cd$ , it follows that every call of ENUM( $u_1, \dots, u_i$ ) leads to at least one recursive call of ENUM( $u_1, \dots, u_i, u_{i+1}$ ). Furthermore, there are at most  $cd$  iterations of the loop in line 9 that do *not* lead to a recursive call. As every test in line 10 can be done in time  $O(c)$ , it follows that the time spans until the first recursive call, between the calls, and after the last call are bounded by  $O(c^2d)$ . As the recursion depth is  $c$ , the overall delay between two output tuples is bounded by  $O(c^3d)$ .  $\square$

By using similar techniques as in [38], we obtain the following improved version of Lemma 12.15 where the delay is independent of the degree bound  $d$ .

**Lemma 12.16.** *There is a dynamic algorithm that receives a number  $c \geq 1$ , a degree bound  $d \geq 2$  and a  $\sigma_c$ -db  $\mathcal{G}_0$  of degree  $\leq d$ , and computes within  $t_p = d^{\text{poly}(c)} \cdot |\mathcal{G}_0|$  preprocessing time a data structure that can be updated in time  $d^{\text{poly}(c)}$  and allows to enumerate the query result  $\varphi_c(\mathcal{G})$  with  $O(c^2)$  delay.*

Before proving Lemma 12.16, let us first point out that Lemma 12.16 in combination with Theorem 12.8 directly improves the delay in Theorem 12.14 from  $d^{2^{O(\|\varphi\|)}}$  to  $O(k^2)$ , immediately leading to the following theorem.

**Theorem 12.17.** *There is a dynamic algorithm that receives a schema  $\sigma$ , a degree bound  $d \geq 2$ , a  $k$ -ary FO+MOD[ $\sigma$ ]-query  $\varphi(\bar{x})$  (for some  $k \in \mathbb{N}$ ) and a  $\sigma$ -db  $D_0$  of degree  $\leq d$ , and computes within  $t_p = f(Q, d) \cdot \|D_0\|$  preprocessing time a data structure that can be updated in time  $t_u = f(Q, d)$  and allows to enumerate  $\varphi(D)$  with  $O(k^2)$  delay.*

The rest of the section is devoted to the proof of Lemma 12.16.

## 12. Answering non-Boolean FO+MOD Queries Under Updates

*Proof of Lemma 12.16.* Consider Algorithm 18, which enumerates  $\varphi_c(\mathcal{G})$  with  $O(c^3d)$  delay. To enumerate the tuples with only  $O(c^2)$  delay, we replace the loop in lines 9–10 by a precomputed “skip” function that allows to iterate through all elements in  $T_{i+1}^{\mathcal{G}} \setminus \bigcup_{j=1}^i N^{\mathcal{G}}(u_j)$  with  $O(c)$  delay.

For every  $i \in [c]$  we store all elements of  $T_i^{\mathcal{G}}$  in a doubly linked list and let **void** be an auxiliary element that appears at the end of the list. We let  $\text{first}_i$  be the first element in the list and  $\text{succ}_i(u)$  the successor of  $u \in T_i^{\mathcal{G}}$ . We denote by  $\leq^i$  the linear order induced by this list. We let  $\tilde{E}^{\mathcal{G}}$  be the symmetric closure of  $E^{\mathcal{G}}$ , i.e.,  $\tilde{E}^{\mathcal{G}} = E^{\mathcal{G}} \cup \{(v, u) : (u, v) \in E^{\mathcal{G}}\}$ . For every  $i \in [c]$  we define the function

$$\text{skip}_i(y, V) := \min \left\{ z \in T_i^{\mathcal{G}} \cup \{\text{void}\} : y \leq^i z \text{ and for all } v \in V, (v, z) \notin \tilde{E}^{\mathcal{G}} \right\},$$

which assigns to every  $V \subseteq \text{adom}(\mathcal{G})$  with  $|V| \leq c-1$ , and every  $y \in T_i^{\mathcal{G}}$  the next node in  $T_i^{\mathcal{G}}$  that is not adjacent to any vertex in  $V$ .

Using these functions, our improved enumeration algorithm is given in Algorithm 19. Below, we show that we can access the values  $\text{skip}_i(y, V)$  in time  $O(c)$ . By the same analysis as given in the proof of Lemma 12.15 it then follows that Algorithm 19 enumerates  $\varphi_c(\mathcal{G})$  with  $O(c^2)$  delay.

---

**Algorithm 19** Enumeration procedure with delay  $O(c^2)$ .

---

```

1: for all  $(u_1, \dots, u_s) \in \mathcal{S}$  do
2:    $\text{ENUM}(u_1, \dots, u_s)$ .
3: Output the end-of-enumeration message EOE.
4:
5: function  $\text{ENUM}(u_1, \dots, u_i)$ 
6:   if  $i = c$  then
7:     output the tuple  $(u_1, \dots, u_c)$ .
8:   else
9:      $y \leftarrow \text{skip}_{i+1}(\text{first}_{i+1}, \{u_1, \dots, u_i\})$ 
10:    while  $y \neq \text{void}$  do
11:       $\text{ENUM}(u_1, \dots, u_i, y)$ .
12:       $y \leftarrow \text{skip}_{i+1}(\text{succ}_{i+1}(y), \{u_1, \dots, u_i\})$ .
```

---

What remains to show is that we can access the values  $\text{skip}_i(y, V)$  for all  $i, y, V$  in time  $O(c)$  and maintain them with  $d^{\text{poly}(c)}$  update time. At first sight, this is not clear at all, because the domain of  $\text{skip}_i$  has size  $\Omega(|\text{adom}(\mathcal{G})|^c)$ . In what follows, we show that for every  $y$ , the number of distinct values that  $\text{skip}_i(y, V)$  can take is bounded by  $d^{\text{poly}(c)}$  and that we can store them in a look-up table with update time  $d^{\text{poly}(c)}$ .

To illustrate the main idea, let us start with a simple example. We want to enumerate  $\varphi_4$  on a coloured graph  $\mathcal{H}$  with four vertex colours blue, red, yellow and green (in this order) and analyse the call of  $\text{ENUM}(b, r, y)$ , which is supposed to enumerate all green nodes  $g_i$  that are not adjacent to any of the nodes  $b, r$  and  $y$ . The relevant part of  $\mathcal{H}$  is depicted in Figure 12.1.

### 12.5. Enumerating Results of FO+MOD Queries Under Updates

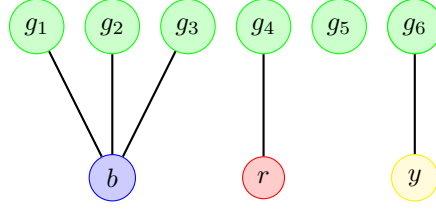


Figure 12.1.: Illustration of the relevant part of graph  $\mathcal{H}$ .

The enumeration procedure starts by considering the first element  $g_1$  in the list of green vertices, but the first element in the actual output is  $g_5 = \text{skip}_4(g_1, \{b, r, y\})$ . Therefore, we have to skip the irrelevant vertices  $g_1, \dots, g_4$ .

To do this, we want to know the neighbours of the vertices that we skip ( $b$  and  $r$  in our example) when looking at  $g_1$ . For this purpose, we define inductively new sorts of edges  $E_4^1 \subseteq E_4^2 \subseteq \dots$  that connect green vertices  $g \in \{g_1, \dots, g_6\}$  with  $\tilde{E}$ -neighbours of skipped vertices. In our example, we first have to skip  $g_1$ , because it is  $\tilde{E}$ -connected to  $b$  and to be able to handle this, we let  $E_4^1$  be the set of tuples  $(g_i, v) \in \tilde{E}^{\mathcal{H}}$  (see Figure 12.2).

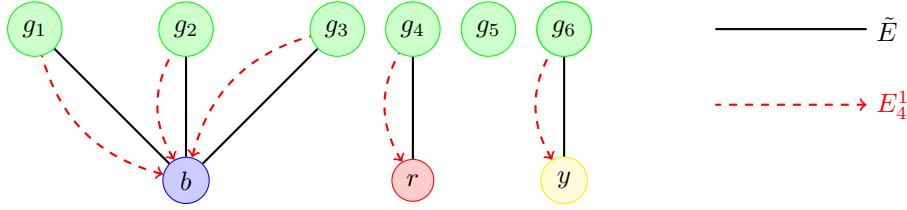


Figure 12.2.:  $\tilde{E}$ -edges and  $E_4^1$ -edges in our example.

After realising that even more vertices ( $g_2$  and  $g_3$ ) are excluded by  $b$ , the next try would be  $g_4$ . However, this vertex is excluded by its  $\tilde{E}$ -neighbour  $r$ , so we have to take  $r$  into account when computing the skip value for  $g_1$  and indicate this by the  $E_4^2$ -edge  $(g_1, r)$  (see Figure 12.3). This immediately leads to an inductive definition:  $E_4^2$  contains all pairs of vertices that are already in  $E_4^1$  or connected by a path as shown in Figure 12.4.

The idea outlined above can be formalised as follows. For  $i, j \in [c]$ , we define inductively the auxiliary edge sets  $E_i^j$ :

$$E_i^1 := \left\{ (y, u) : y \in C_i^{\mathcal{G}} \text{ and } (y, u) \in \tilde{E}^{\mathcal{G}} \right\} \quad \text{and}$$

$$E_i^{j+1} := E_i^j \cup \left\{ (y, u) : \begin{array}{l} \text{there are } v, z \text{ with } (y, v) \in E_i^j, \\ (v, z) \in \tilde{E}^{\mathcal{G}}, \text{ and } (\text{succ}_i(z), u) \in \tilde{E}^{\mathcal{G}} \end{array} \right\}$$

Now we define for every  $y \in C_i^{\mathcal{G}}$  the set

$$S_i^y := \{ u : (y, u) \in E_i^c \}.$$

## 12. Answering non-Boolean FO+MOD Queries Under Updates

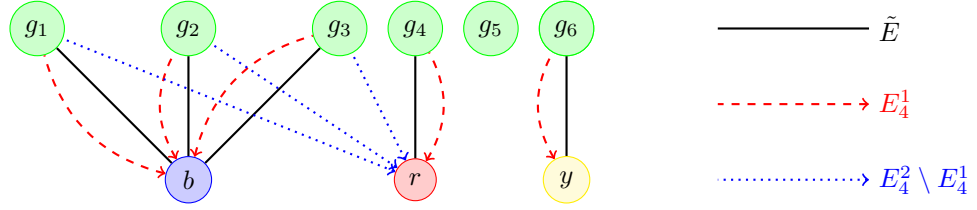


Figure 12.3.:  $\tilde{E}$ -edges,  $E_4^1$ -edges and  $E_4^2$ -edges in our example.

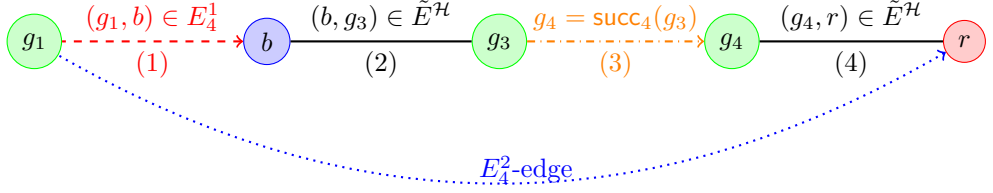


Figure 12.4.: Introducing an  $E_4^2$ -edge between  $g_1$  and  $r$ .

Note that  $|S_i^y| = O(d^{2c})$ . The following claim states that the elements of  $S_i^y$  are the only ones we need to take into account when computing  $\text{skip}_i(y, V)$ .

**Claim 12.18.** *For all  $i \leq c$ ,  $y \in C_i^G \cup \{\text{void}\}$  and  $V \subseteq \text{adom}(\mathcal{G})$  with  $|V| \leq c-1$  it holds that*

$$\text{skip}_i(y, V) = \text{skip}_i(y, V \cap S_i^y). \quad (12.10)$$

*Proof.* The proof is identical to the proof of Claim 1 in [38]. For the reader's convenience, the proof is included here. If  $c = 1$  or  $y = \text{void}$ , the claim is trivial. Hence assume that  $c \geq 2$ ,  $y \neq \text{void}$  and let  $z := \text{skip}_i(y, V \cap S_i^y)$ . By definition we have  $y \leq^i z \leq^i \text{skip}_i(y, V)$  and therefore we have to show  $z \geq^i \text{skip}_i(y, V)$ , which holds if and only if  $(u, z) \notin \tilde{E}^G$  for all  $u \in V \setminus S_i^y$ . If  $z = y$ , the claim clearly holds as all  $\tilde{E}^G$ -neighbours of  $y$  are contained in  $S_i^y$ . Hence we have  $z >^i y$  and let  $z' \geq^i y$  be the predecessor of  $z$ , i.e.,  $z = \text{succ}_i(z')$ . Now assume for contradiction that there is an  $u \in V \setminus S_i^y$  such that  $(*) (u, z) \in \tilde{E}^G$ . Note that since  $z' <^i z = \text{skip}_i(y, V \cap S_i^y)$ , there is a  $v \in V \cap S_i^y$  such that  $(**) (v, z') \in \tilde{E}^G$ . In the following we show that  $(***) (y, v) \in E_i^{c-1}$ . Note that this finishes the proof of the claim, as by the definition of  $E_i^c$ , the statements  $(*)$ ,  $(**)$  and  $(***)$  imply that  $u \in S_i^y$ , contradicting the assumption that  $u \in V \setminus S_i^y$ .

To show that  $(y, v) \in E_i^{c-1}$ , let

$$V_j := \left\{ v' \in V : (y, v') \in E_i^j \right\} \quad (12.11)$$

for all  $j \in [c]$ . Note that  $V_c = V \cap S_i^y$ . Furthermore, if there is a  $j < c$  with  $V_j = V_{j+1}$ ,



## 12.5. Enumerating Results of FO+MOD Queries Under Updates

then we have

$$V_j = V_{j+1} = \dots = V_c = V \cap S_i^y. \quad (12.12)$$

Since  $|V| \leq c-1$  and  $u \in V \setminus S_i^y$ , we have  $|V \cap S_i^y| \leq c-2$ . In particular, it holds that  $V_{c-1} = V \cap S_i^y$ . Since  $v \in V \cap S_i^y$ , it holds that  $v \in V_{c-1}$  and thus  $(y, v) \in E_i^{c-1}$ .  $\square$

In our dynamic algorithm we maintain an array that allows random access to the values  $\text{skip}_i(y, S')$  for all  $y \in C_i^G$  and all  $S' \subseteq S_i^y$  of size at most  $c-1$ . By Claim 12.18 we can then compute  $\text{skip}_i(y, V)$  by first computing  $S' = V \cap S_i^y$  and then looking up  $\text{skip}_i(y, S')$ . This can be done in time  $O(c)$ . The next claim states that we can efficiently maintain the sets  $S_i^y$ .

**Claim 12.19.** *There is a data structure that*

1. *stores the elements from the sets  $S_i^y$  and all subsets  $S' \subseteq S_i^y$  of cardinality at most  $c-1$ ,*
2. *allows to test membership in these sets in time  $O(1)$  and*
3. *can be updated in time  $d^{\text{poly}(c)}$  after every update of the form insert  $C_i(v)$ , delete  $C_i(v)$ , insert  $E(u, v)$ , and delete  $E(u, v)$ .*

*Proof.* Note that  $u \in S_i^y \iff (y, u) \in E_i^c$ . We store the edge sets  $E_i^c$  for all  $i \in [c]$  in adjacency lists and additionally maintain arrays to allow constant-time access to all list entries. In particular, the adjacency list contains for all  $y \in V$  a list of elements from  $S_i^y$  and access the elements in  $S_i^y$  in constant time. Moreover, as the size of  $S_i^y$  is bounded by  $O(d^{2c})$ , the number of subsets  $S' \subseteq S_i^y$  of cardinality at most  $c-1$  is bounded by  $O(d^{2c^2})$ . Consequently, we can provide constant-time access to all these subsets  $S'$ .

On every insertion or deletion of an edge in  $E^G$ , as well as every insertion or deletion of a vertex in  $C_i^G$ , we update the edge sets  $E_i^c$  according to their inductive definition (see Algorithm 20 for a function computing the  $E_i^c$ -edges). To do this efficiently, we use a breadth-first search starting from  $u$  and  $v$ , for every tuple  $(u, v)$  that has changed in relation  $E^G$  (or for every  $u$  that has changed the membership in  $C_i^G$ ), up to depth  $3c$  to identify the relevant nodes that are affected by the change. For all such nodes  $y$  we first delete the  $E_i^c$ -edges where  $y$  is the first component and then we insert the edges  $E_i^c$  that are enumerated by  $\text{COMPUTE}_E(y, i, c)$  in Algorithm 20. It is straightforward to see that the tuples yield in  $\text{COMPUTE}_E(y, i, c)$  belongs to  $E_i^c$  and that this can be done in time  $d^{\text{poly}(c)}$  as the degree of the edge sets is bounded by  $d^{\text{poly}(c)}$ . For every vertex  $y$  where the set  $S_i^y$  is recomputed we compute the set  $S' \subseteq S_i^y$  of cardinality at most  $c-1$ . This can be done in time  $d^{\text{poly}(c)}$ .  $\square$

In our data structure we store the values  $\text{skip}_i(y, S')$  for every  $i \in [c]$ ,  $y \in T_i^G$  and for all sets  $S' \subseteq S_i^y$  of cardinality at most  $c-1$ . On every insertion or deletion of an edge, we update the sets  $S_i^y$  and their subsets  $S'$  of cardinality at most  $c-1$  and update affected values of  $\text{skip}_i(y, S')$ . According to Claim 12.19 this can be done in time  $d^{\text{poly}(c)}$ .

---

**Algorithm 20** Compute the sets  $E_i^j$ .

---

```

1: function COMPUTEE( $y, i, j$ )
2:   if  $j = 1$  then
3:     if  $y \in C_i^{\mathcal{G}}$  then
4:       for  $(y, u) \in \tilde{E}^{\mathcal{G}}$  do
5:         yield  $(y, u)$ 
6:   else  $\triangleright i > 1$ 
7:     for  $(y, v)$  in COMPUTEE( $y, i, j - 1$ ) do
8:       yield  $(y, v)$   $\triangleright$  This tuple belongs to  $E_i^{j-1}$ 
9:       for  $(v, z) \in \tilde{E}^{\mathcal{G}}$  do
10:        for  $(\text{succ}_i(z), u) \in \tilde{E}^{\mathcal{G}}$  do
11:          yield  $(y, u)$ 

```

---

We do the same on updates of the form `insert`  $C_i(v)$  and `delete`  $C_i(v)$ , but have to do some additional work, as  $v$  might occur in the image of skip-functions. Upon `insert`  $C_i(v)$ , we insert  $v$  at the beginning of the list  $C_i$ . This ensures that existing skip values will not be affected. Afterwards, we compute the set  $S_i^v$  and the values  $\text{skip}_i(v, S')$  for all  $S' \subseteq S_i^v$  of cardinality at most  $c-1$ . Again, this can be done in time  $d^{\text{poly}(c)}$ .

If we receive the update `delete`  $C_i(v)$ , then we have to recompute all skip values  $\text{skip}_i(y, S')$  that point to  $v$ . Note that (since  $\mathcal{G}$  has degree  $\leq d$ ) this is only the case for nodes  $y \leq^i v$  whose distance from  $v$  w.r.t.  $\text{succ}_i$  is at most  $(c-1)d$ . Hence, it suffices to recompute  $\text{skip}_i(y, S')$  for at most  $(c-1)d$  vertices  $y$  and all  $S' \subseteq S_i^y$  of cardinality at most  $c-1$ . This can be done in time  $d^{\text{poly}(c)}$ . By Claim 12.18, all this suffices to access the value for  $\text{skip}_i(y, V)$  in time  $O(c)$ . This concludes the proof of Lemma 12.16.  $\square$

## 12.6. Refining the enumeration routine

In the previous section we have presented a dynamic algorithm that allows to enumerate with delay  $O(k^2)$  the result  $Q(D)$  of a  $k$ -ary FO+MOD-query  $Q(x_1, \dots, x_k)$  on a database  $D$  of degree at most  $d$  (see Theorem 12.17). In this section, we generalise this to provide the following functionality. Upon input of a tuple  $\bar{a} = (a_1, \dots, a_\ell) \in \text{dom}^\ell$  for some  $\ell \in [k]$ , we would like to be able to enumerate all tuples  $\bar{b} = (b_1, \dots, b_k)$  in  $Q(D)$  whose first  $\ell$  components coincide with  $\bar{a}$ . In fact, since we already know that the first  $\ell$  components of  $\bar{b}$  coincide with  $\bar{a}$ , we only output the remaining components  $(b_{\ell+1}, \dots, b_k)$  of  $\bar{b}$ .

Upon input of the tuple  $\bar{a}$ , our algorithm spends some *preparation time*  $t_{\text{prep}}$ , after which it outputs all the desired result tuples with delay  $t_d$ .

For formulating this chapter's main result, the following notation will be convenient. Given  $k \geq 1$  and  $\ell \in [k]$ , we let  $m := k - \ell$ . For a tuple  $\bar{x} = (x_1, \dots, x_k)$  of  $k$  pairwise distinct variables we let  $\bar{z} = (z_1, \dots, z_\ell) := (x_1, \dots, x_\ell)$  and  $\bar{y} = (y_1, \dots, y_m) := (x_{\ell+1}, \dots, x_k)$ .

## 12.6. Refining the enumeration routine

For a  $k$ -ary  $\text{FO}+\text{MOD}[\sigma]$ -query  $Q(\bar{x}) = Q(\bar{z}, \bar{y})$  and a tuple  $\bar{a} = (a_1, \dots, a_\ell) \in \mathbf{dom}^\ell$  we let  $Q(\bar{a}, \bar{y})$  be the  $m$ -ary query which, when evaluated on a  $\sigma$ -db  $D$  returns the result set

$$\{ (b_1, \dots, b_m) : (a_1, \dots, a_\ell, b_1, \dots, b_m) \in Q(D) \}.$$

This section is devoted to the proof of the following theorem.

**Theorem 12.20.** *There is a dynamic algorithm that receives a schema  $\sigma$ , a degree bound  $d \geq 2$ , a  $k$ -ary  $\text{FO}+\text{MOD}[\sigma]$ -query  $Q(\bar{x})$  (for some  $k \in \mathbb{N}_{\geq 1}$ ), a number  $\ell \in [k]$  and a  $\sigma$ -db  $D_0$  of degree  $\leq d$ , and computes within  $t_p = f(Q, d) \cdot \|D_0\|$  preprocessing time a data structure that can be updated in time  $t_u = f(Q, d)$  and provides the following functionality: Upon input of an arbitrary tuple  $\bar{a} = (a_1, \dots, a_\ell) \in \mathbf{dom}^\ell$ , after  $t_{\text{prep}} = f(Q, d)$  preparation time it enumerates with delay  $O(k^2)$  all result tuples of  $Q(\bar{a}, \bar{y})$  on  $D$ .*

For proving Theorem 12.20, we will use the following variant of Lemma 12.16, which deals with the queries

$$\varphi_J(z_{j_1}, \dots, z_{j_c}) := \bigwedge_{j \in J} T_j(z_j) \wedge \bigwedge_{j, j' \in J, j \neq j'} \neg E(z_j, z_{j'}) \quad (12.13)$$

for all sets  $J = \{j_1, \dots, j_c\} \subseteq [2c]$  of cardinality  $|J| = c$ .

**Lemma 12.21.** *There is a dynamic algorithm that receives a number  $c \geq 1$ , a degree bound  $d \geq 2$  and a  $\sigma_{2c}$ -db  $\mathcal{G}_0$  of degree  $\leq d$ , and computes within  $t_p = d^{\text{poly}(c)} \cdot |\mathcal{G}_0|$  preprocessing time a data structure that can be updated in time  $d^{\text{poly}(c)}$  and provides the following functionality. Upon input of an arbitrary set  $J \subseteq [2c]$  of size  $c$ , after  $d^{\text{poly}(c)}$  preparation time it enumerates  $\varphi_J(\mathcal{G})$  with delay  $O(c^2)$ .*

*Proof.* We use the dynamic algorithm provided by the proof of Lemma 12.16 for input  $2c$  instead of  $c$ . When receiving an update command, we apply this algorithm's **update** routine.

Recall that the dynamic data structure constructed in the proof of Lemma 12.16 (for  $2c$  instead of  $c$ ) stores the set  $I \subseteq [2c]$  of *small* colour classes in  $\mathcal{G}$ , i.e.,  $i \in I$  iff  $|C_i^\mathcal{G}| \leq 2cd$ .

Upon input of a set  $J$  we proceed as follows. Compute the set  $I' := I \cap J$  of all small colour classes that are relevant for the query  $\varphi_J$ . For simplicity, let us assume that  $I' = \{1, \dots, s\}$  for some  $s \leq c$  and  $J \setminus I' = \{s+1, \dots, c\}$  (otherwise, we rename the colours accordingly). Within preparation time  $d^{\text{poly}(c)}$  we can compute the set

$$\mathcal{S} := \left\{ (u_1, \dots, u_s) \in C_1^\mathcal{G} \times \dots \times C_s^\mathcal{G} : \begin{array}{l} (u_j, u_{j'}) \notin E^\mathcal{G}, \\ \text{for all } j, j' \in [s] \text{ with } j \neq j' \end{array} \right\}.$$

To enumerate the query result  $\varphi_J(\mathcal{G})$ , we then use Algorithm 19 for this set  $\mathcal{S}$  (without any changes; in particular, in lines 6 and 7 we don't replace  $c$  by  $2c$ , but keep the value  $c$ ).

Revisiting the proof of Lemma 12.16 it is straightforward to verify that the resulting dynamic data structure provides the desired functionality within the claimed preparation time and delay.  $\square$

## 12. Answering non-Boolean FO+MOD Queries Under Updates

Based on the above lemma, we can prove the following variant of Lemma 12.9.

**Lemma 12.22.** *For  $d', c \in \mathbb{N}$  let  $t_i(c, d') := (d')^{\text{poly}(c)}$ ,  $t_u(c, d') := (d')^{\text{poly}(c)}$  and  $t_d(c, d') := O(c^2)$ .*

*For every schema  $\sigma$  and every  $d \geq 2$  the following holds. Let  $r \geq 0$ ,  $k \geq 1$  and fix  $d' := d^{2k^2(2r+1)}$  and  $\tilde{t}_x := \max_{c=1}^k t_x(c, d')$  for  $t_x \in \{t_i, t_u, t_c, t_d\}$ .*

(1) *Let  $\tau$  be a  $d$ -bounded  $r$ -type with  $k$  centres and let  $\ell \in [k]$ . There is a dynamic algorithm which within initialisation time  $\tilde{t}_{\text{init}}$  builds a data structure that can be updated in time at most  $\tilde{t}_{\text{update}} d^{O(k^2 r + k \|\sigma\|)} + 2^{O(\|\sigma\| k^2 d^{2r+2})}$  and provides the following functionality for  $\sigma$ -dbs of degree at most  $d$ . Upon input of an arbitrary tuple  $\bar{a} \in \mathbf{dom}^\ell$ , after  $\tilde{t}_{\text{update}} 2^{O(\|\sigma\| k^2 d^{2r+2})}$  preparation time it enumerates with delay  $O(\tilde{t}_{\text{delay}} + k)$  all result tuples of  $\text{sph}_\tau(\bar{a}, \bar{y})$  on  $D$ .*

(2) *Let  $s \geq 0$  and let  $\chi_1, \dots, \chi_s$  be arbitrary sentences of schema  $\sigma$ , and assume we have available for each  $j \in [s]$  a dynamic algorithm with initialisation time  $t'_i$  and update time  $t'_u$  that allows to check within answer time  $t'_{\text{ans}}$  whether or not  $D \models \chi_j$  for  $d$ -bounded  $\sigma$ -dbs  $D$ .*

*Let  $\bar{x} = (x_1, \dots, x_k)$  be a tuple of pairwise distinct variables and let  $\psi(\bar{x})$  be a Boolean combination of the sentences  $\chi_1, \dots, \chi_s$  and of  $d$ -bounded sphere-formulas of radius at most  $r$  (over  $\sigma$ ).*

*There is a dynamic algorithm which within initialisation time*

$$s(t'_i + t'_{\text{ans}}) + 2^{(kd^{r+1})^{O(\|\sigma\|)}} (\text{poly}(\|\psi\|) + \tilde{t}_{\text{init}}),$$

*builds a data structure that can be updated in time*

$$s(t'_u + t'_{\text{ans}}) + 2^{(k^2 d^{2r+2})^{O(\|\sigma\|)}} (\text{poly}(\|\psi\|) + \tilde{t}_{\text{update}} d^{O(k^2 r + k \|\sigma\|)})$$

*and provides the following functionality. Upon input of a tuple  $\bar{a} \in \mathbf{dom}^\ell$ , after  $2^{(k^2 d^{2r+2})^{O(\|\sigma\|)}} \cdot \tilde{t}_{\text{update}}$  preparation time it enumerates with delay  $O(\tilde{t}_{\text{delay}} + k)$  all result tuples of  $\psi(\bar{a}, \bar{y})$  on  $D$ .*

*Proof.* For the proof of Part (1) we proceed in a similar way as in the proof of Lemma 12.9(1) and utilise the dynamic algorithm provided by Lemma 12.21.

Upon input of a tuple  $\bar{a} = (a_1, \dots, a_\ell) \in \mathbf{dom}^\ell$  we want to enumerate all tuples  $\bar{b} \in \mathbf{dom}^m$  such that  $D \models \text{sph}_\tau(\bar{a}, \bar{b})$ . We use the same notation as in the proof of Lemma 12.9(1). In particular, recall that we consider the formula

$$\text{conn-sph}_\tau(\bar{x}) := \bigwedge_{j \in [c]} \text{sph}_{\tau_j, \nu_j}(\bar{x}_j) \wedge \bigwedge_{j \neq j'} \neg \text{dist}_{\leq 2r+1}^{k_j, k_{j'}}(\bar{x}_j, \bar{x}_{j'})$$

and the coloured graph  $\mathcal{G}_D$  of degree at most  $d'$  with

$$\begin{aligned} T_j^{\mathcal{G}_D} &:= \{ v_{\bar{b}} : \bar{b} \in \text{adom}(D)^{k_j} \text{ with } (\mathcal{N}_r^D(\bar{b}), \bar{b}) \cong \tau_j \}, \quad \text{for all } j \in [c], \text{ and} \\ E^{\mathcal{G}_D} &:= \left\{ (v_{\bar{b}}, v_{\bar{b}'} ) \in V^2 : \text{dist}^D(\bar{b}, \bar{b}') \leq 2r+1 \right\}, \end{aligned}$$

## 12.6. Refining the enumeration routine

where  $V := \bigcup_{j \in [c]} T_j^{\mathcal{G}_D}$  and recall that  $\tau_j := \tau_{j, \nu_j}$  is a connected  $r$ -type for each  $j \in [c]$ .

Recall that  $(x_1, \dots, x_\ell) = (z_1, \dots, z_\ell) = \bar{z}$  are the variables the input tuple  $\bar{a}$  will be assigned to. W.l.o.g. we assume that there is a number  $\kappa \in \{1, \dots, c\}$  such that the following is true. For each  $j \in [\kappa]$ , the tuple  $\bar{x}_j$  contains at least one of the variables in  $\bar{z}$  and for each  $j \in [c] \setminus [\kappa]$ , the tuple  $\bar{x}_j$  contains none of the variables in  $\bar{z}$ , i.e., it only consists of variables in  $\bar{y} = (y_1, \dots, y_m) = (x_{\ell+1}, \dots, x_k)$ .

For an input tuple  $\bar{a}$ , we extend  $\mathcal{G}_D$  by  $c$  further colours  $C_{c+1}, \dots, C_{2c}$  to obtain the following  $\sigma_{2c}$ -db  $\mathcal{G}_{D, \bar{a}}$ . The edge relation  $E$  and the colours  $C_1, \dots, C_c$  of  $\mathcal{G}_{D, \bar{a}}$  are the same as those of  $\mathcal{G}_D$ . For each  $j \in \{1, \dots, \kappa\}$  we let

$$C_{c+j}^{\mathcal{G}_{D, \bar{a}}} := \left\{ v_{\bar{b}} \in C_j^{\mathcal{G}_D} : \begin{array}{l} \text{for every position } \nu \text{ in } \bar{x}_j \text{ which consists of a vari-} \\ \text{able } z_i \text{ in } \bar{z}, \text{ the entry of } \bar{b} \text{ at position } \nu \text{ is } a_i \end{array} \right\},$$

and for each  $j \in \{\kappa+1, \dots, c\}$  we let  $C_{c+j}^{\mathcal{G}_{D, \bar{a}}} := \emptyset$ .

To enumerate the result tuples of  $\text{sph}_\tau(\bar{a}, \bar{y})$  we use the **enumerate** routine provided by Lemma 12.21 upon input

$$J := \{c+1, \dots, c+\kappa\} \cup \{\kappa+1, \dots, c\}$$

to enumerate the set  $\varphi_J(\mathcal{G}_{D, \bar{a}})$ . Note that according to the definition of the sets  $C_{c+j}^{\mathcal{G}_{D, \bar{a}}}$ , the set  $\varphi_J(\mathcal{G}_{D, \bar{a}})$  contains exactly those tuples in  $(v_{\bar{b}_1}, \dots, v_{\bar{b}_c}) \in \varphi_J(\mathcal{G}_D)$  where for every position  $\nu$  in some  $\bar{x}_j$  that consists of a variable  $z_i$  in  $\bar{z}$ , the entry of  $\bar{b}_j$  at position  $\nu$  is  $a_i$ . Therefore, after suitably re-ordering the components of the tuples  $(\bar{b}_1, \dots, \bar{b}_c)$  and dropping the components that correspond to  $\bar{a}$ , we obtain the desired result tuples for  $\text{sph}_\tau(\bar{a}, \bar{y})$  on  $D$ .

It remains to show how to construct and maintain  $\mathcal{G}_{D, \bar{a}}$  when the database is updated or when a new tuple  $\bar{a}$  is given as input.

The **update** routine proceeds in exactly the same way as in the proof of Lemma 12.9(1), uses the **update** routine provided by Lemma 12.21, but does not bother to update the colours  $C_{c+1}, \dots, C_{2c}$ .

Upon input of a tuple  $\bar{a}$ , we use the preparation time to first delete all elements from the colours  $C_{c+1}, \dots, C_{2c}$ , and then insert all the elements that belong to the sets  $C_{c+1}^{\mathcal{G}_{D, \bar{a}}}, \dots, C_{c+\kappa}^{\mathcal{G}_{D, \bar{a}}}$ . The details can be carried out as follows.

Note that according to the definition of  $C_{c+j}^{\mathcal{G}_{D, \bar{a}}}$ , the following is true for each  $j \in [\kappa]$ : For every element  $v_{\bar{b}}$  in  $C_{c+j}^{\mathcal{G}_{D, \bar{a}}}$ , some component of the tuple  $\bar{b}$  is equal to a component  $a_i$  of the tuple  $\bar{a}$ . Since  $(\mathcal{N}_r^D(\bar{b}), \bar{b}) \cong \tau_j$  and  $\tau_j$  is a connected  $r$ -type, every component of the tuple  $\bar{b}$  belongs to  $\mathcal{N}_{r'}^D(a_i)$  for  $r' := r + (k-1)(2r+1)$ . In particular, from Lemma 10.2 we obtain that  $|C_{c+j}^{\mathcal{G}_{D, \bar{a}}}| \leq |\mathcal{N}_{r'}^D(a_i)|^k \leq d^{k^2(2r+1)}$ . Thus, the first step of deleting all elements that are still present in  $C_{c+j}$  (as an artifact of a previous enumeration request) can be accomplished in time  $d^{O(k^2r)} \tilde{t}_{\text{update}}$ .

Afterwards, we proceed as follows to insert into  $C_{c+j}$  all elements that do belong to  $C_{c+j}^{\mathcal{G}_{D, \bar{a}}}$ . We compute the set  $U := \mathcal{N}_{r'}^D(a_i)$  and test for every tuple  $\bar{b} \in U^{|\bar{x}_j|}$  whether the following two conditions hold:

## 12. Answering non-Boolean FO+MOD Queries Under Updates

- (1) The  $r'$ -type  $(\mathcal{N}_{r'}^D(\bar{b}), \bar{b})$  is isomorphic to  $\tau_j$ , and
- (2) for every position  $\nu$  in  $\bar{x}_j$  which consists of a variable  $z_\mu$  in  $\bar{z}$ , the entry of  $\bar{b}$  at position  $\nu$  is  $a_\mu$ .

If both conditions are met, we insert the vertex  $v_{\bar{b}}$  into  $C_{c+j}$  by using the **update** routine provided by Lemma 12.21. Note that this constructs the correct set  $C_{c+j}^{\mathcal{G}_D, \bar{\pi}}$ .

To analyse the time needed for this construction, note that by Lemma 10.2,  $|U| \leq d^{r'+1} \leq d^{k(2r+1)} \leq d^{O(kr)}$ . Furthermore,  $U$  can be computed in time  $(d^{r'+1})^{O(\|\sigma\|)} \leq d^{O(kr\|\sigma\|)}$ . The number of tuples  $\bar{b}$  that we consider is at most  $|U|^k \leq d^{O(k^2r)}$ . To check if the first condition for  $\bar{b}$  is met, we use Lemma 10.2(e) to check in time  $2^{O(\|\sigma\|k^2d^{2r+2})}$  whether the  $r$ -type of  $\bar{b}$  is isomorphic to  $\tau_j$ . The second condition can be checked in time  $O(k)$ . In summary, we compute the sets  $C_{c+j}^{\mathcal{G}_D, \bar{\pi}}$  for all  $j \in [\kappa]$  in time  $d^{O(k^2r\|\sigma\|)} \cdot 2^{O(\|\sigma\|k^2d^{2r+2})} \cdot \tilde{t}_{\text{update}} \leq 2^{O(\|\sigma\|k^2d^{2r+2})} \cdot \tilde{t}_{\text{update}}$ .

This completes the proof of part (1).

Let us now turn to the proof of part (2). The proof can be taken verbatim from the proof of Lemma 12.9(2), with the following changes. Instead of using the dynamic algorithms for  $\text{sph}_{\tau_i}(\bar{x})$  provided from Lemma 12.9(1), we now use the dynamic algorithms for  $\text{sph}_{\tau_i}(\bar{z}, \bar{y})$  provided by Lemma 12.22(1). Upon input of a tuple  $\bar{a} \in \text{dom}^\ell$ , we run the preparation phase of the dynamic algorithms for  $\text{sph}_{\tau_i}(\bar{a}, \bar{y})$  for all  $i \in I$ , and afterwards, we loop through all  $i \in I$  and enumerate the result tuples of  $\text{sph}_{\tau_i}(\bar{a}, \bar{y})$  on  $D$ . This completes the proof of Lemma 12.22.  $\square$

By using Lemma 12.22(2) instead of Lemma 12.9(2), we obtain the following analogue to Theorem 12.8.

**Lemma 12.23.** *For  $d', c \in \mathbb{N}$  let  $t_i(c, d') = (d')^{\text{poly}(c)}$ ,  $t_u(c, d') = (d')^{\text{poly}(c)}$ , and  $t_d(c, d') = O(c^2)$ . There is a dynamic algorithm that receives a schema  $\sigma$ , a degree bound  $d \geq 2$ , a  $k$ -ary FO+MOD[ $\sigma$ ]-query  $Q(\bar{x})$  (for some  $k \in \mathbb{N}_{\geq 1}$ ), and a number  $\ell \in [k]$ . Within initialisation time  $\hat{t}_{\text{init}} \cdot f(Q, d)$ , this algorithm builds a data structure that can be updated in time  $\hat{t}_{\text{update}} \cdot f(Q, d)$  and provides the following functionality for  $\sigma$ -dbs of degree at most  $d$ . Upon input of an arbitrary tuple  $\bar{a} \in \text{dom}^\ell$ , after  $\hat{t}_{\text{update}} \cdot f(Q, d)$  preparation time it enumerates with delay  $O(\hat{t}_{\text{delay}} + k)$  all result tuples of  $Q(\bar{a}, \bar{y})$  on  $D$ , where  $\hat{t}_x = \max_{c=1}^k t_x(c, d^{2^{O(\|\sigma\|)}})$  for  $t_x \in \{t_i, t_u, t_d\}$ .*

*Proof.* The proof can be taken verbatim from the proof of Theorem 12.8, with the following change. Instead of applying part (2) of Lemma 12.9, we now use the dynamic algorithm provided by part (2) of Lemma 12.22.  $\square$

*Proof of Theorem 12.20.*

The result is an immediate consequence of Lemma 12.23.  $\square$

## 12.7. Enumerating the Difference

In this chapter we consider the update routine called **update\_and\_report\_diff** which, immediately after performing the database update, reports the difference between the new query result and the old query result, i.e., it first enumerates all tuples in  $Q(D^+) \setminus Q(D^-)$  (terminated by the end-of-enumeration message **EOE**) and then all tuples in  $Q(D^-) \setminus Q(D^+)$  (again, terminated by the message **EOE**), where  $D^-$  and  $D^+$  denote the database before and after performing the given update command. This chapter's main result reads as follows.

**Theorem 12.24.** *There is a dynamic algorithm that receives a schema  $\sigma$ , a degree bound  $d \geq 2$ , a  $k$ -ary  $\text{FO}+\text{MOD}[\sigma]$ -query  $Q(\bar{y})$  (for some  $k \in \mathbb{N}$ ) and a  $\sigma$ -db  $D_0$  of degree  $\leq d$ , and computes within  $t_p = f(Q, d) \cdot \|D_0\|$  preprocessing time a data structure providing an **update\_and\_report\_diff** routine which, upon input of a command  $\text{update } R(\bar{a})$  for  $\text{update} \in \{\text{insert}, \text{delete}\}$ ,  $R \in \sigma$  and  $\bar{a} \in \text{dom}^{\text{ar}(R)}$ , updates the data structure within time  $t_u = f(Q, d)$  and then enumerates the difference between the new and the old query result with delay  $O(k^2)$ .*

*Proof.* We use Theorem 12.20 for the following queries. For each  $R \in \sigma$  and each  $\text{update} \in \{\text{insert}, \text{delete}\}$  we let  $\ell := \text{ar}(R)$  and fix a tuple  $\bar{z} = (z_1, \dots, z_\ell)$  of pairwise distinct variables that do not occur in  $Q(\bar{x})$ . We construct  $(\ell+k)$ -ary  $\text{FO}+\text{MOD}[\sigma]$ -queries

$$Q_{\text{update } R}^+(\bar{z}, \bar{x}) \quad \text{and} \quad Q_{\text{update } R}^-(\bar{z}, \bar{x})$$

such that the following is true for all  $\sigma$ -dbs  $D$  and all  $\bar{a} \in \text{dom}^\ell$ : if  $D^+$  is the  $\sigma$ -db obtained from  $D^- \neq D^+$  by performing the update operation  $\text{update } R(\bar{a})$ , then

- the set of result tuples of  $Q_{\text{update } R}^+(\bar{a}, \bar{x})$  on  $D^+$  is exactly the set  $Q(D^+) \setminus Q(D^-)$ , and
- the set of result tuples of  $Q_{\text{update } R}^-(\bar{a}, \bar{x})$  on  $D^+$  is exactly the set  $Q(D^-) \setminus Q(D^+)$ .

Before constructing these queries let us explain how they can be used to finish the proof of Theorem 12.24. Let  $\Psi$  be the set consisting of the queries  $Q_{\text{update } R}^+$  and  $Q_{\text{update } R}^-$  for  $\text{update} \in \{\text{insert}, \text{delete}\}$  and  $R \in \sigma$ . We use in parallel for each  $\psi$  in  $\Psi$  the dynamic algorithm provided by Theorem 12.20. Upon input of an update operation  $\text{update } R(\bar{a})$ , the **update\_and\_report\_diff** routine proceeds as follows. Let  $D = D^-$  be the database before executing the update command.

In case that the given update command does not change the database (i.e., the operation intends to delete a tuple that does not belong to the database relation  $R^D$  or it intends to insert a tuple that already belongs to  $R^D$  or it intends to insert a tuple that would result in a database that exceeds the given degree bound  $d$ ), then all the data structures remain unchanged and the routine just outputs “EOE”.

Otherwise, we proceed as follows. First, consider each query  $\psi$  in  $\Psi$  and perform the **update** routine upon input “ $\text{update } R(\bar{a})$ ” of the dynamic algorithm provided by Theorem 12.20 for the query  $\psi$ . Let  $D^+$  be the updated database. We then use the functionality provided by Theorem 12.20 to perform in parallel the preparation phase

## 12. Answering non-Boolean FO+MOD Queries Under Updates

for the queries  $Q_{\text{update } R}^+(\bar{z}, \bar{x})$  and  $Q_{\text{update } R}^-(\bar{z}, \bar{x})$  upon input of the tuple  $\bar{a}$ . Afterwards, we first enumerate the result tuples of  $Q_{\text{update } R}^+(\bar{a}, \bar{x})$  on  $D^+$  and then enumerate the result tuples of  $Q_{\text{update } R}^-(\bar{a}, \bar{x})$  on  $D^+$ .

All that remains to be done to finish the proof of Theorem 12.24 is to construct the queries  $Q_{\text{update } R}^+(\bar{z}, \bar{x})$  and  $Q_{\text{update } R}^-(\bar{z}, \bar{x})$ . The idea is straightforward: we let

$$\begin{aligned} Q_{\text{update } R}^+(\bar{z}, \bar{x}) &:= Q(\bar{x}) \wedge \neg Q_{\text{update } R}^{\text{old}}(\bar{z}, \bar{x}) \quad \text{and} \\ Q_{\text{update } R}^-(\bar{z}, \bar{x}) &:= Q_{\text{update } R}^{\text{old}}(\bar{z}, \bar{x}) \wedge \neg Q(\bar{x}), \end{aligned}$$

where  $Q_{\text{update } R}^{\text{old}}(\bar{z}, \bar{x})$  is a formula for which the following is true: if  $D^+ \neq D^-$  is obtained from  $D^-$  by performing the update command  $\text{update } R(\bar{a})$ , then evaluating  $Q_{\text{update } R}^{\text{old}}(\bar{a}, \bar{x})$  on  $D^+$  simulates the evaluation of  $Q(\bar{x})$  on  $D^-$ . A somewhat annoying technical detail in the construction of these formulas is that the semantics of quantifiers takes into account the database's active domain, and the database's active domain might be changed by the update command.

Let us first consider the (easier) case that  $\text{update} = \text{insert}$ . We let

$$\begin{aligned} \alpha(\bar{z}, y) &:= \exists y_1 \cdots \exists y_\ell \left( \bigvee_{i=1}^{\ell} y=y_i \wedge R(y_1, \dots, y_\ell) \wedge \neg \bigwedge_{i=1}^{\ell} y_i=z_i \right) \vee \\ &\quad \bigvee_{S \in \sigma \setminus \{R\}} \exists y_1 \cdots \exists y_{\text{ar}(S)} \left( \bigvee_{i=1}^{\text{ar}(S)} y=y_i \wedge S(y_1, \dots, y_{\text{ar}(S)}) \right). \end{aligned}$$

If  $D^+ \neq D^-$  is obtained from  $D^-$  by performing the update command  $\text{insert } R(\bar{a})$ , then the result set of  $\alpha(\bar{a}, y)$  on  $D^+$  is exactly the active domain of  $D^-$ . Therefore, we can choose  $Q_{\text{insert } R}^{\text{old}}(\bar{z}, \bar{x})$  to be the query obtained from the input query  $Q(\bar{x})$  by replacing every atomic subformula of the form  $R(u_1, \dots, u_\ell)$  with the formula  $(R(u_1, \dots, u_\ell) \wedge \neg \bigwedge_{i=1}^{\ell} u_i=z_i)$  and by relativising every quantification to a variable  $y$  to those  $y$  that satisfy  $\alpha(\bar{z}, y)$ , i.e., we replace every subformula of the form  $\exists y \vartheta$  (or  $\exists^{i \bmod m} y \vartheta$ ) with the formula  $\exists y (\alpha(\bar{z}, y) \wedge \vartheta)$  (or  $\exists^{i \bmod m} y (\alpha(\bar{z}, y) \wedge \vartheta)$ ). It is straightforward to verify that the resulting formula  $Q_{\text{insert } R}^{\text{old}}(\bar{z}, \bar{x})$  expresses the desired property.

Let us now turn to the case where  $\text{update} = \text{delete}$ . The problem here is that  $\text{adom}(D^-)$  contains all the elements in  $\bar{a} = (a_1, \dots, a_\ell)$ , while some (or, all) of these elements might be missing in  $\text{adom}(D^+)$ , and due to the active domain semantics of FO+MOD, there is no explicit means of enabling quantifiers to range over elements that do not belong to the active domain. To overcome this, let  $J \subseteq [\ell]$  be a set of indices such that  $|J| = |\{a_1, \dots, a_\ell\}|$  and  $\{a_1, \dots, a_\ell\} = \{a_j : j \in J\}$ . By induction on the construction of formulas we define for every FO+MOD[ $\sigma$ ]-query  $\vartheta(\bar{y})$  that does not contain any of the variables in  $\bar{z} = (z_1, \dots, z_\ell)$  an FO+MOD[ $\sigma$ ]-query  $\hat{\vartheta}_J(\bar{z}, \bar{y})$  such that the set of result tuples of  $\hat{\vartheta}_J(\bar{a}, \bar{y})$  on  $D^+$  is exactly the set  $\vartheta(D^-)$ . To achieve this, we proceed as follows:

- if  $\vartheta$  is of the form  $R(y_1, \dots, y_\ell)$ , then  $\hat{\vartheta}_J := (R(y_1, \dots, y_\ell) \vee \bigwedge_{i=1}^{\ell} y_i=z_i)$



## 12.7. Enumerating the Difference

- if  $\vartheta$  is of the form  $y_1=y_2$  or of the form  $S(y_1, \dots, y_{\text{ar}(S)})$  with  $S \in \sigma \setminus \{R\}$ , then  $\hat{\vartheta}_J := \vartheta$
- if  $\vartheta$  is of the form  $\neg\theta$ , then  $\hat{\vartheta}_J := \neg\hat{\theta}_J$
- if  $\vartheta$  is of the form  $(\theta' \vee \theta'')$ , then  $\hat{\vartheta}_J := (\hat{\theta}'_J \vee \hat{\theta}''_J)$
- if  $\vartheta$  is of the form  $\exists y \theta$ , then  $\hat{\vartheta}_J := (\vartheta \vee \bigvee_{j \in J} \hat{\theta}_J \frac{z_j}{y})$ , where  $\hat{\theta}_J \frac{z_j}{y}$  is the formula obtained from  $\hat{\theta}_J$  by replacing every free occurrence of the variable  $y$  by the variable  $z_j$
- if  $\vartheta$  is of the form  $\exists^{i \bmod m} y \theta$ , then  $\hat{\vartheta}_J := \bigvee_{(i_1, i_2) \in I} \xi_{(i_1, i_2)}$ , where  $I$  is the set of all  $(i_1, i_2) \in \{0, \dots, m-1\}^2$  with  $i_1 + i_2 \equiv i \bmod m$ , and

$$\xi_{(i_1, i_2)} := \exists^{i_1 \bmod m} y \left( \hat{\theta}_J \wedge \bigwedge_{j \in J} \neg y = z_j \right) \wedge \bigvee_{\substack{J' \subseteq J \\ |J'| = i_2}} \left( \bigwedge_{j \in J'} \hat{\theta}_J \frac{z_j}{y} \wedge \bigwedge_{j \in J \setminus J'} \neg \hat{\theta}_J \frac{z_j}{y} \right).$$

It is straightforward to verify that the query  $\hat{\vartheta}_J(\bar{z}, \bar{y})$  indeed has the desired meaning. Thus, we can choose

$$Q_{\text{delete } R}^{\text{old}}(\bar{z}, \bar{x}) := \bigvee_{J \subseteq [\ell]} \left( \hat{Q}_J \wedge \bigwedge_{\substack{i, j \in J \\ i \neq j}} \neg z_i = z_j \wedge \bigwedge_{i \in [\ell] \setminus J} \bigvee_{j \in J} z_i = z_j \right).$$

This completes the proof of Theorem 12.24.  $\square$



## 13. Conclusion

The main results in the first part show that we can maintain the results of q-hierarchical queries under updates, and they can be tested and counted in constant time and enumerated with constant delay, after a linear time preprocessing, they can be updated in constant time. Furthermore, it was shown that for the testing problem, there is a data structure that can be maintained with constant update time for the more expressive language of t-hierarchical conjunctive queries. Moreover, we showed that if we spend more update time, i.e., logarithmic time in the size of the database, we can also support the  $j$ th routine. The results for conjunctive queries have been lifted to so called q-hierarchical (t-hierarchical) unions of conjunctive queries, for which we can support the enumeration and the testing problem under updates and moreover, we considered a subset of q-hierarchical unions of conjunctive queries, the strongly exhaustively q-hierarchical union of conjunctive queries for which we can support the  $j$ th problem under updates.

In the papers [16] and [21] it is shown that there are lower bounds for conjunctive queries in the task of dynamic query evaluation. If a self-join free conjunctive query is not q-hierarchical then there is no algorithm with arbitrary preprocessing time and update time sublinear in the size of the databases's active domain, that enumerates sublinear delay, unless the OMv-conjecture (a conjecture on the hardness of online matrix-vector problem [57]) fails. Similar results exist for answering a Boolean conjunctive query and for the counting problem for non-Boolean conjunctive queries. Moreover, there is still no characterisation of the conjunctive queries with aggregates that cannot be maintained under updates. Furthermore, for the  $j$ th problem there is no lower bound for unions of conjunctive queries.

Another task we considered is that the preparation of learning a polynomial regression function can be done in constant update time if the training data is taken from a query result of a q-hierarchical query. An obvious future task is to investigate if there are other machine learning models where the training data can be taken from the query, where the data structure is maintained under updates.

In this thesis we considered running times with respect to data complexity. A further future task is to figure out the fine grained complexity of answering conjunctive queries under updates in the combined complexity, i.e., additionally in the size of the query. Moreover, another interesting task is to give an amortized analysis on how the delay depends on the query.

The main results in the second part show that in the dynamic setting (i.e., allowing database updates), the results of  $k$ -ary FO+MOD-queries on bounded degree databases can be tested and counted in constant time and enumerated with constant delay, after linear time preprocessing and with constant update time. Here, "constant time" refers to data complexity and is of size  $\text{poly}(k)$  concerning the delay and the time for testing

### 13. Conclusion

and counting. The time for performing a database update is 3-fold exponential in the size of the query and the degree bound, and is worst-case optimal.

The starting point of our algorithms is to decompose the given query into a query in Hanf normal form, using a recent result of [55]. This normal form is only available for the setting with a fixed maximum degree bound  $d$ , i.e., the setting considered in this part.

Recently, Kuske and Schweikardt [70] introduced a new kind of Hanf normal form for a variant of *first-order logic with counting* that contains and extends Libkin's logic  $\text{FO}(\text{Cnt})$  [72] and Grohe's logic  $\text{FO}+\text{C}$  [48]. As an application it is shown in [70] that the techniques presented in Part II (Chapter 10 - 12) can be lifted from  $\text{FO}+\text{MOD}$  to first-order logic with counting terms and numerical predicates  $\text{FOC}(\mathbb{P})$ .

An obvious future task is to investigate to which extent further query evaluation results that are known for the static setting can be lifted to the dynamic setting. More specifically: are there efficient dynamic algorithms for evaluating (i.e., answering, testing, counting or enumerating) results of first-order queries on other sparse classes of databases (e.g. planar, bounded treewidth, bounded expansion, nowhere dense) or databases of low degree, lifting the “static” results accumulated in [61, 49, 38, 93] to the dynamic setting?

# Bibliography

- [1] Serge Abiteboul, Marcelo Arenas, Pablo Barceló, Meghyn Bienvenu, Diego Calvanese, Claire David, Richard Hull, Eyke Hüllermeier, Benny Kimelfeld, Leonid Libkin, Wim Martens, Tova Milo, Filip Murlak, Frank Neven, Magdalena Ortiz, Thomas Schwentick, Julia Stoyanovich, Jianwen Su, Dan Suciu, Victor Vianu, and Ke Yi. Research directions for principles of data management (dagstuhl perspectives workshop 16151). *Dagstuhl Manifestos*, 7(1):1–29, 2018.
- [2] Serge Abiteboul, Richard Hull, and Victor Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [3] Georgii M. Adelson-Velskii and Evgenii M. Landis. An algorithm for the organization of information (in russian). In *Doklady Akademii Nauk SSSR*, pages 146:263–266, 1962.
- [4] Antoine Amarilli, Pierre Bourhis, Stefan Mengel, and Matthias Niewerth. Enumeration on trees with tractable combined complexity and efficient updates. In *Proceedings of the 38th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2019, Amsterdam, The Netherlands, June 30 - July 5, 2019.*, pages 89–103, 2019.
- [5] Molham Aref, Balder ten Cate, Todd J. Green, Benny Kimelfeld, Dan Olteanu, Emir Pasalic, Todd L. Veldhuizen, and Geoffrey Washburn. Design and implementation of the logicblox system. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, Melbourne, Victoria, Australia, May 31 - June 4, 2015*, pages 1371–1382, 2015.
- [6] Albert Atserias, Martin Grohe, and Dániel Marx. Size bounds and query plans for relational joins. *SIAM J. Comput.*, 42(4):1737–1767, 2013.
- [7] Guillaume Bagan. *Algorithmes et complexité des problèmes d’énumération pour l’évaluation de requêtes logiques. (Algorithms and complexity of enumeration problems for the evaluation of logical queries)*. PhD thesis, University of Caen Normandy, France, 2009.
- [8] Guillaume Bagan, Arnaud, and Etienne Grandjean. On acyclic conjunctive queries and constant delay enumeration. In *Proceedings of the 16th Annual Conference of the EACSL, CSL’07, Lausanne, Switzerland, September 11–15, 2007*, pages 208–222, 2007.
- [9] Guillaume Bagan, Arnaud Durand, Etienne Grandjean, and Frédéric Olive. Computing the jth solution of a first-order query. *ITA*, 42(1):147–164, 2008.

## BIBLIOGRAPHY

- [10] Nurzhan Bakibayev, Tomáš Kociský, Dan Olteanu, and Jakub Zavodny. Aggregation and ordering in factorised databases. *PVLDB*, 6(14):1990–2001, 2013.
- [11] Andrey Balmin, Yannis Papakonstantinou, and Victor Vianu. Incremental validation of XML documents. *ACM Trans. Database Syst.*, 29(4):710–751, 2004.
- [12] Daniel Barbará and Tomasz Imielinski. Sleepers and workaholics: Caching strategies in mobile environments. *VLDB J.*, 4(4):567–602, 1995.
- [13] Michael Benedikt and H. Jerome Keisler. Expressive power of unary counters. In *Database Theory - ICDT '97, 6th International Conference, Delphi, Greece, January 8-10, 1997, Proceedings*, pages 291–305, 1997.
- [14] Michael Benedikt and Leonid Libkin. Exact and approximate aggregation in constraint query. In *Proceedings of the Eighteenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, May 31 - June 2, 1999, Philadelphia, Pennsylvania, USA*, pages 102–113, 1999.
- [15] Christoph Berkholz, Jens Keppeler, and Nicole Schweikardt. Answering conjunctive queries under updates. *CoRR*, abs/1702.06370, 2017.
- [16] Christoph Berkholz, Jens Keppeler, and Nicole Schweikardt. Answering conjunctive queries under updates. In *Proceedings of the 36th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS'17, Chicago, IL, USA, May 14-19, 2017*, pages 303–318, 2017. Full version available at <http://arxiv.org/abs/1702.06370>.
- [17] Christoph Berkholz, Jens Keppeler, and Nicole Schweikardt. Answering FO+MOD queries under updates on bounded degree databases. *CoRR*, abs/1702.08764, 2017.
- [18] Christoph Berkholz, Jens Keppeler, and Nicole Schweikardt. Answering FO+MOD queries under updates on bounded degree databases. In Michael Benedikt and Giorgio Orsi, editors, *20th International Conference on Database Theory, ICDT 2017, March 21-24, 2017, Venice, Italy*, volume 68 of *LIPIcs*, pages 8:1–8:18. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2017.
- [19] Christoph Berkholz, Jens Keppeler, and Nicole Schweikardt. Answering ucqs under updates and in the presence of integrity constraints. *CoRR*, abs/1709.10039, 2017.
- [20] Christoph Berkholz, Jens Keppeler, and Nicole Schweikardt. Answering FO+MOD queries under updates on bounded degree databases. *ACM Trans. Database Syst.*, 43(2):7:1–7:32, 2018.
- [21] Christoph Berkholz, Jens Keppeler, and Nicole Schweikardt. Answering ucqs under updates and in the presence of integrity constraints. In *21st International Conference on Database Theory, ICDT 2018, March 26-29, 2018, Vienna, Austria*, pages 8:1–8:19, 2018.

## BIBLIOGRAPHY

- [22] Henrik Björklund, Wouter Gelade, and Wim Martens. Incremental xpath evaluation. *ACM Trans. Database Syst.*, 35(4):29:1–29:43, 2010.
- [23] Johann Brault-Baron. *De la pertinence de l'énumération : complexité en logiques propositionnelle et du premier ordre. (The relevance of the list: propositional logic and complexity of the first order)*. PhD thesis, University of Caen Normandy, France, 2013.
- [24] Andrei A. Bulatov, Víctor Dalmau, Martin Grohe, and Dániel Marx. Enumerating homomorphisms. *J. Comput. Syst. Sci.*, 78(2):638–650, 2012.
- [25] Peter J Cameron. *Combinatorics: topics, techniques, algorithms*. Cambridge University Press, 1994.
- [26] Nofar Carmeli and Markus Kröll. On the enumeration complexity of unions of conjunctive queries. In *Proceedings of the 38th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2019, Amsterdam, The Netherlands, June 30 - July 5, 2019.*, pages 134–148, 2019.
- [27] Surajit Chaudhuri, Ravi Krishnamurthy, Spyros Potamianos, and Kyuseok Shim. Optimizing queries with materialized views. In *Proceedings of the Eleventh International Conference on Data Engineering, March 6-10, 1995, Taipei, Taiwan*, pages 190–200, 1995.
- [28] Hubie Chen and Stefan Mengel. A trichotomy in the complexity of counting answers to conjunctive queries. In *Proceedings of the 18th International Conference on Database Theory, ICDT'15, Brussels, Belgium, March 23–27, 2015*, pages 110–126, 2015.
- [29] Hubie Chen and Stefan Mengel. Counting answers to existential positive queries: A complexity classification. In *Proceedings of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS'16, San Francisco, CA, USA, June 26 – July 01, 2016*, pages 315–326, 2016.
- [30] Rada Chirkova and Jun Yang. Materialized views. *Foundations and Trends in Databases*, 4(4):295–405, 2012.
- [31] Sara Cohen. Containment of aggregate queries. *SIGMOD Record*, 34(1):77–85, 2005.
- [32] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms (3. ed.)*. MIT Press, 2009.
- [33] Víctor Dalmau and Peter Jonsson. The complexity of counting homomorphisms seen from the other side. *Theoretical Computer Science, Volume 329, Issues 1–3*, pages 315–323, December 2004.

## BIBLIOGRAPHY

- [34] Nilesh Dalvi and Dan Suciu. The dichotomy of conjunctive queries on probabilistic structures. In *Proceedings of the 26th ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, PODS'07, Beijing, China, June 11–13, 2007*, pages 293–302. ACM, 2007.
- [35] Alin Dobra, Minos N. Garofalakis, Johannes Gehrke, and Rajeev Rastogi. Processing complex aggregate queries over data streams. In *Proceedings of the 2002 ACM SIGMOD International Conference on Management of Data, Madison, Wisconsin, USA, June 3–6, 2002*, pages 61–72, 2002.
- [36] Arnaud Durand and Etienne Grandjean. First-order queries on structures of bounded degree are computable with constant delay. *ACM Trans. Comput. Log.*, 8(4), 2007.
- [37] Arnaud Durand and Stefan Mengel. Structural tractability of counting of solutions to conjunctive queries. *Theory of Computing Systems, Volume 57, Issue 4*, pages 1202–1249, November 2015. Full version of an ICDT 2013 paper.
- [38] Arnaud Durand, Nicole Schweikardt, and Luc Segoufin. Enumerating answers to first-order queries over databases of low degree. In *Proceedings of the 33rd ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS'14, Snowbird, UT, USA, June 22–27, 2014*, pages 121–131, 2014.
- [39] Arnaud Durand and Yann Strozecki. Enumeration complexity of logical query problems with second-order variables. In *Computer Science Logic, 25th International Workshop / 20th Annual Conference of the EACSL, CSL 2011, September 12–15, 2011, Bergen, Norway, Proceedings*, pages 189–202, 2011.
- [40] Kousha Etessami and Neil Immerman. Tree canonization and transitive closure. *Inf. Comput.*, 157(1-2):2–24, 2000.
- [41] Jon Feldman, S. Muthukrishnan, Anastasios Sidiropoulos, Clifford Stein, and Zoya Svitkina. On the complexity of processing massive, unordered, distributed data. *CoRR*, abs/cs/0611108, 2006.
- [42] Robert Fink and Dan Olteanu. Dichotomies for queries with negation in probabilistic databases. *ACM Trans. Database Syst.*, 41(1):4:1–4:47, 2016.
- [43] Peter Flach. *Machine learning: the art and science of algorithms that make sense of data*. Cambridge University Press, 2012.
- [44] Markus Frick and Martin Grohe. The complexity of first-order and monadic second-order logic revisited. *Ann. Pure Appl. Logic*, 130(1–3):3–31, 2004.
- [45] Gianluigi Greco and Francesco Scarcello. Structural tractability of enumerating CSP solutions. *Constraints*, 18(1):38–74, 2013.



- [46] Gianluigi Greco and Francesco Scarcello. Counting solutions to conjunctive queries: structural and hybrid tractability. In *Proceedings of the 33rd ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS'14, Snowbird, UT, USA, June 22–27, 2014*, pages 132–143, 2014.
- [47] Martin Grohe. The complexity of homomorphism and constraint satisfaction problems seen from the other side. *Journal of the ACM (JACM)*, Volume 54, Issue 1, Article No. 1, March 2007.
- [48] Martin Grohe. *Descriptive Complexity, Canonisation, and Definable Graph Structure Theory*. Lecture Notes in Logic. Association for Symbolic Logic in conjunction with Cambridge University Press, to appear. Preliminary version available at <https://www.lii.rwth-aachen.de/de/13-mitarbeiter/professoren/39-book-descriptive-complexity.html>.
- [49] Martin Grohe, Stephan Kreutzer, and Sebastian Siebertz. Deciding first-order properties of nowhere dense graphs. *J. ACM*, 64(3):17:1–17:32, 2017. Conference version: in *Proceedings of the 46th ACM Symposium on Theory of Computing (STOC'14)*, pp.89–98, 2014.
- [50] Martin Grohe and Nicole Schweikardt. First-order query evaluation with cardinality conditions. In *Proceedings of the 37th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS'18, Houston, TX, USA, June 10–15, 2018*, 2018.
- [51] Martin Grohe, Thomas Schwentick, and Luc Segoufin. When is the evaluation of conjunctive queries tractable? In *Proceedings on 33rd Annual ACM Symposium on Theory of Computing, STOC'01, Heraklion, Crete, Greece, July 6–8, 2001*, pages 657–666, 2001.
- [52] Ashish Gupta, Venky Harinarayan, and Dallan Quass. Aggregate-query processing in data warehousing environments. In *VLDB'95, Proceedings of 21th International Conference on Very Large Data Bases, September 11–15, 1995, Zurich, Switzerland.*, pages 358–369, 1995.
- [53] Ashish Gupta, Inderpal Singh Mumick, and V. S. Subrahmanian. Maintaining views incrementally. In *Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data, SIGMOD'93, Washington, D.C., USA, May 25–28, 1993*, pages 157–166. ACM, 1993.
- [54] Isabelle Guyon, Steve Gunn, Masoud Nikraves, and Lofti A Zadeh. *Feature extraction: foundations and applications*, volume 207. Springer, 2008.
- [55] Lucas Heimberg, Dietrich Kuske, and Nicole Schweikardt. Hanf normal form for first-order logic with unary counting quantifiers. In *Proceedings of the 31st Annual ACM/IEEE Symposium on Logic in Computer Science, LICS '16, New York, NY, USA, July 5–8, 2016*, pages 277–286, 2016.

## BIBLIOGRAPHY

- [56] Lauri Hella, Leonid Libkin, Juha Nurmonen, and Limsoon Wong. Logics with aggregate operators. *J. ACM*, 48(4):880–907, 2001.
- [57] Monika Henzinger, Sebastian Krinninger, Danupon Nanongkai, and Thatchaphol Saranurak. Unifying and strengthening hardness for dynamic problems via the online matrix-vector multiplication conjecture. In *Proceedings of the 47th Annual ACM on Symposium on Theory of Computing, STOC’15, Portland, OR, USA, June 14–17, 2015*, pages 21–30, 2015.
- [58] Botong Huang, Matthias Boehm, Yuanyuan Tian, Berthold Reinwald, Shirish Tatikonda, and Frederick R. Reiss. Resource elasticity for large-scale machine learning. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, Melbourne, Victoria, Australia, May 31 - June 4, 2015*, pages 137–152, 2015.
- [59] Muhammad Idris, Martín Ugarte, and Stijn Vansummeren. The dynamic Yannakakis algorithm: Compact and efficient query processing under updates. In *Proceedings of the 2017 ACM International Conference on Management of Data, SIGMOD’17, Chicago, IL, USA, May 14–19, 2017*, pages 1259–1274, 2017.
- [60] Wojciech Kazana and Luc Segoufin. First-order query evaluation on structures of bounded degree. *Logical Methods in Computer Science*, 7(2), 2011.
- [61] Wojciech Kazana and Luc Segoufin. Enumeration of first-order queries on classes of structures with bounded expansion. In *Proceedings of the 32nd ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS 2013, New York, NY, USA, June 22–27, 2013*, pages 297–308, 2013.
- [62] Mahmoud Abo Khamis, Hung Q. Ngo, XuanLong Nguyen, Dan Olteanu, and Maximilian Schleich. In-database learning with sparse tensors. In *Proceedings of the 37th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, Houston, TX, USA, June 10–15, 2018*, pages 325–340, 2018.
- [63] Mahmoud Abo Khamis, Hung Q. Ngo, Christopher Ré, and Atri Rudra. Joins via geometric resolutions: Worst-case and beyond. In *Proceedings of the 34th ACM Symposium on Principles of Database Systems, PODS’15, Melbourne, Victoria, Australia, May 31 - June 4, 2015*, pages 213–228, 2015.
- [64] Benny Kimelfeld and Christopher Ré. A relational framework for classifier engineering. In *Proceedings of the 36th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2017, Chicago, IL, USA, May 14–19, 2017*, pages 5–20, 2017.
- [65] Donald Ervin Knuth. *The art of computer programming*, volume 3. Pearson Education, 1997.
- [66] Christoph Koch. Incremental query evaluation in a ring of databases. In *Proceedings of the 29th ACM SIGMOD-SIGACT-SIGART Symposium on Principles*

## BIBLIOGRAPHY

- of Database Systems, PODS'10, June 6–11, 2010, Indianapolis, Indiana, USA, pages 87–98, 2010.
- [67] Christoph Koch, Yanif Ahmad, Oliver Kennedy, Milos Nikolic, Andres Nötzli, Daniel Lupei, and Amir Shaikhha. Dbtoaster: higher-order delta processing for dynamic, frequently fresh views. *VLDB J.*, 23(2):253–278, 2014.
  - [68] Christoph Koch, Daniel Lupei, and Val Tannen. Incremental view maintenance for collection programming. In *Proceedings of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS'16, San Francisco, CA, USA, June 26 – July 01, 2016*, pages 75–90, 2016.
  - [69] Paraschos Koutris and Dan Suciu. Parallel evaluation of conjunctive queries. In *Proceedings of the 30th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS 2011, June 12-16, 2011, Athens, Greece*, pages 223–234, 2011.
  - [70] Dietrich Kuske and Nicole Schweikardt. First-order logic with counting: At least, weak Hanf normal forms always exist and can be computed! In *Proceedings of the 32nd Annual ACM/IEEE Symposium on Logic in Computer Science, LICS'17*, 2017. Full version available at CoRR abs/1703.01122.
  - [71] Alon Y. Levy, Divesh Srivastava, and Thomas Kirk. Data model and query evaluation in global information systems. *J. Intell. Inf. Syst.*, 5(2):121–143, 1995.
  - [72] Leonid Libkin. *Elements of Finite Model Theory*. Texts in Theoretical Computer Science. An EATCS Series. Springer, 2004.
  - [73] Leonid Libkin and Limsoon Wong. New techniques for studying set languages, bag languages and aggregate functions. In *Proceedings of the Thirteenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, May 24-26, 1994, Minneapolis, Minnesota, USA*, pages 155–166, 1994.
  - [74] Katja Losemann and Wim Martens. MSO queries on trees: enumerating answers under updates. In *Joint Meeting of the Twenty-Third EACSL Annual Conference on Computer Science Logic (CSL) and the Twenty-Ninth Annual ACM/IEEE Symposium on Logic in Computer Science (LICS), CSL-LICS '14, Vienna, Austria, July 14 - 18, 2014*, pages 67:1–67:10, 2014.
  - [75] Eugene M. Luks. Isomorphism of graphs of bounded valence can be tested in polynomial time. *J. Comput. Syst. Sci.*, 25(1):42–65, 1982.
  - [76] Dániel Marx. Tractable hypergraph properties for constraint satisfaction and conjunctive queries. *Journal of the ACM (JACM)*, Volume 60, Issue 6, Article No. 42, November 2013.
  - [77] Bernard M. E. Moret and Henry D. Shapiro. *Algorithms from P to NP: Volume 1: Design & Efficiency*. Benjamin-Cummings, 1991.

## BIBLIOGRAPHY

- [78] Hung Q. Ngo, Dung T. Nguyen, Christopher Re, and Atri Rudra. Beyond worst-case analysis for joins with minesweeper. In *Proceedings of the 33rd ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS'14, Snowbird, UT, USA, June 22–27, 2014*, pages 234–245, 2014.
- [79] Hung Q. Ngo, Ely Porat, Christopher Ré, and Atri Rudra. Worst-case optimal join algorithms: [extended abstract]. In *Proceedings of the 31st ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS'12, Scottsdale, AZ, USA, May 20–24, 2012*, pages 37–48, 2012.
- [80] Matthias Niewerth and Luc Segoufin. Enumeration of MSO queries on strings with constant delay and logarithmic updates. In *Proceedings of the 37th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, Houston, TX, USA, June 10–15, 2018*, pages 179–191, 2018.
- [81] Milos Nikolic, Mohammad Dashti, and Christoph Koch. How to win a hot dog eating contest: Distributed incremental view maintenance with batch updates. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD'16, San Francisco, CA, USA, June 26 – July 01, 2016*, pages 511–526, 2016.
- [82] Dan Olteanu and Jakub Závodný. Size bounds for factorised representations of query results. *ACM Trans. Database Syst.*, 40(1):2:1–2:44, 2015.
- [83] Sushant Patnaik and Neil Immerman. Dyn-FO: A parallel, dynamic complexity class. *J. Comput. Syst. Sci.*, 55(2):199–209, 1997.
- [84] Christopher Ré, Divy Agrawal, Magdalena Balazinska, Michael J. Cafarella, Michael I. Jordan, Tim Kraska, and Raghu Ramakrishnan. Machine learning and databases: The sound of things to come or a cacophony of hype? In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, Melbourne, Victoria, Australia, May 31 - June 4, 2015*, pages 283–284, 2015.
- [85] Sylvio Rüdian. So funktioniert die seo-ki, 2019. (In German) Available at <https://seor1d.com/seo-ki>, Accessed: 2019-09-18.
- [86] Maximilian Schleich, Dan Olteanu, and Radu Ciucanu. Learning linear regression models over factorized joins. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, pages 3–18, 2016.
- [87] Nicole Schweikardt, Luc Segoufin, and Alexandre Vigny. Enumeration for FO queries over nowhere dense graphs. In *Proceedings of the 37th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS'18, Houston, TX, USA, June 10–15, 2018*, 2018.
- [88] Thomas Schwentick and Thomas Zeume. Dynamic complexity: recent updates. *SIGLOG News*, 3(2):30–52, 2016.

## BIBLIOGRAPHY

- [89] Detlef Seese. Linear time computable problems and first-order descriptions. *Mathematical Structures in Computer Science*, 6(6):505–526, 1996.
- [90] Luc Segoufin. Enumerating with constant delay the answers to a query. In *Proceedings of the 16th International Conference on Database Theory, ICDT'13, Genoa, Italy, March 18–22, 2013*, pages 10–20, 2013.
- [91] Luc Segoufin. A glimpse on constant delay enumeration (invited talk). In *Proceedings of the 31st International Symposium on Theoretical Aspects of Computer Science, STACS'14, March 5–8, 2014, Lyon, France*, pages 13–27, 2014.
- [92] Luc Segoufin. Constant delay enumeration for conjunctive queries. *SIGMOD Record*, 44(1):10–17, 2015.
- [93] Luc Segoufin and Alexandre Vigny. Constant delay enumeration for FO queries over databases with local bounded expansion. In *20th International Conference on Database Theory, ICDT 2017, March 21–24, 2017, Venice, Italy*, pages 20:1–20:16, 2017.
- [94] Divesh Srivastava, Shaul Dar, H. V. Jagadish, and Alon Y. Levy. Answering queries with aggregation using views. In *VLDB'96, Proceedings of 22th International Conference on Very Large Data Bases, September 3–6, 1996, Mumbai (Bombay), India*, pages 318–329, 1996.
- [95] Todd L. Veldhuizen. Triejoin: A simple, worst-case optimal join algorithm. In *Proceedings of the 17th International Conference on Database Theory, ICDT'14, Athens, Greece, March 24–28, 2014*, pages 96–106, 2014.
- [96] Thomas Zeume and Thomas Schwentick. Dynamic conjunctive queries. In *Proceedings of the 17th International Conference on Database Theory, ICDT'14, Athens, Greece, March 24–28, 2014*, pages 38–49, 2014.



# Acknowledgements

Firstly, I would like to express my sincere gratitude to my advisor Prof. Dr. Nicole Schweikardt for the continuous support of my Ph.D study, her patience, motivation and offering me a place in her group. Her guidance helped me in all the time of research and writing of this thesis.

I thank Prof. Dr. Stefan Kratsch and Prof. Dr. Wim Martens very much for serving as a reviewer of this thesis.

Of course, I want to thank to all my former and current colleagues in the research group “Logic in Computer Science” for their support and the agreeable atmosphere. In particular, I want to thank Gesine Pergl and Eva Sandig for the administrative support and Petra Kämpfer for the technical support.

Moreover, I thank Philipp Badenhopp, Benedikt Konstantin Gräßle, Michael Rücker, Sylvio Rüdian, Benjamin Schlotter, Dr. Markus Schmid, Dr. Sandra Schulz and Viktor Well for reading parts of this thesis, their helpful comments and their support.





# Sebstständigkeitserklärung

Hiermit erkläre ich, die Dissertation selbstständig und nur unter Verwendung der angegebenen Hilfen und Hilfsmittel angefertigt zu haben. Ich habe mich nicht anderwärts um einen Doktorgrad in dem Promotionsfach beworben und besitze keinen entsprechenden Doktorgrad.

Die Promotionsordnung der Mathematisch-Naturwissenschaftlichen Fakultät, veröffentlicht im Amtlichen Mitteilungsblatt der Humboldt-Universität zu Berlin Nr. 42 am 11. Juli 2018, habe ich zur Kenntnis genommen.

Berlin, den 12. Dezember 2019